



THE UNIVERSITY *of* EDINBURGH

## Edinburgh Research Explorer

### Büchi Objectives in Countable MDPs

**Citation for published version:**

Kiefer, S, Mayr, R, Shirmohammadi, M & Totzke, P 2019, Büchi Objectives in Countable MDPs. in C Baier, I Chatzigiannakis, P Flocchini & S Leonardi (eds), *46th International Colloquium on Automata, Languages, and Programming (ICALP 2019)*. vol. 132, 119, LIPICS, vol. 132, Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, Germany, pp. 119:1–119:14, 46th International Colloquium on Automata, Languages and Programming, Patras, Greece, 8/07/19. <https://doi.org/10.4230/LIPIcs.ICALP.2019.119>

**Digital Object Identifier (DOI):**

[10.4230/LIPIcs.ICALP.2019.119](https://doi.org/10.4230/LIPIcs.ICALP.2019.119)

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Peer reviewed version

**Published In:**

46th International Colloquium on Automata, Languages, and Programming (ICALP 2019)

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



# Büchi Objectives in Countable MDPs

**Stefan Kiefer**

University of Oxford, UK

**Richard Mayr**

University of Edinburgh, UK

**Mahsa Shirmohammadi**

CNRS & IRIF, FR

**Patrick Totzke**

University of Liverpool, UK

---

## Abstract

We study countably infinite Markov decision processes with Büchi objectives, which ask to visit a given subset  $F$  of states infinitely often. A question left open by T.P. Hill in 1979 [10] is whether there always exist  $\varepsilon$ -optimal Markov strategies, i.e., strategies that base decisions only on the current state and the number of steps taken so far. We provide a negative answer to this question by constructing a non-trivial counterexample. On the other hand, we show that Markov strategies with only 1 bit of extra memory are sufficient.

**2012 ACM Subject Classification** Theory of computation  $\rightarrow$  Random walks and Markov chains; Mathematics of computing  $\rightarrow$  Probability and statistics

**Keywords and phrases** Markov decision processes

**Funding** *Stefan Kiefer*: Supported by a Royal Society University Research Fellowship.

*Richard Mayr*: Supported by EPSRC grant EP/M027651/1.

*Mahsa Shirmohammadi*: Supported by PEPS JCJC grant AAPS.

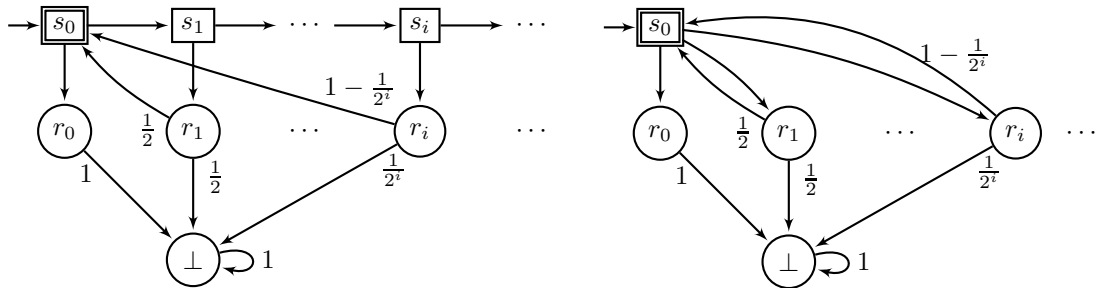
**Acknowledgements** The authors thank anonymous reviewers for their helpful comments.

## 1 Introduction

**Background.** Markov decision processes (MDPs) are a standard model for dynamic systems that exhibit both stochastic and controlled behavior [15]. MDPs play a prominent role in numerous domains, including artificial intelligence and machine learning [18, 17], control theory [4, 1], operations research and finance [5, 16], and formal verification [9, 2]. In an MDP, the system starts in the initial state and makes a sequence of transitions between states. Depending on the type of the current state, either the controller gets to choose an enabled transition (or a distribution over transitions), or the next transition is chosen randomly according to a defined distribution. By fixing a strategy for the controller, one obtains a probability space of runs of the MDP. The goal of the controller is to optimize the expected value of some objective function on the runs.

The type of strategy needed for an optimal (resp.  $\varepsilon$ -optimal) strategy for some objective is also called the *strategy complexity* of the objective. There are different types of strategies, depending on whether one can take the whole history of the run into account (history-dependent; (H)), or whether one is limited to a finite amount of memory (finite memory; (F)) or whether decisions are based only on the current state (memoryless; (M)). Moreover, the strategy type depends on whether the controller can randomize (R) or is limited to deterministic choices (D). The simplest type MD refers to memoryless deterministic strategies. *Markov strategies* are strategies that base their decisions only on the current state and the number of steps in the history or the run. Thus they do use infinite memory, but only in a very restricted form by maintaining an unbounded step-counter. For finite MDPs, there exist optimal MD-strategies for many (but not all) objectives [6, 7, 8, 15], but the picture is more complex for countably infinite MDPs [12, 14, 15].

We study here so-called *Goal* objectives defined via a subset of goal states  $F$ : In the basic Goal objective (also called the *Reachability* objective) one simply wants to reach the set  $F$ . In the *Büchi* objective one wants to visit the set  $F$  infinitely often. For finite MDPs there exist optimal MD-strategies for both these objectives [7, 15]. For countably infinite MDPs, optimal strategies (where they exist) and  $\varepsilon$ -optimal strategies for Reachability can be chosen MD [14, 15]. Similarly, optimal strategies for Büchi (where they exist) can be chosen MD [12]. However,  $\varepsilon$ -optimal strategies for Büchi require infinite memory (cannot be chosen FR); cf. [12, 13].

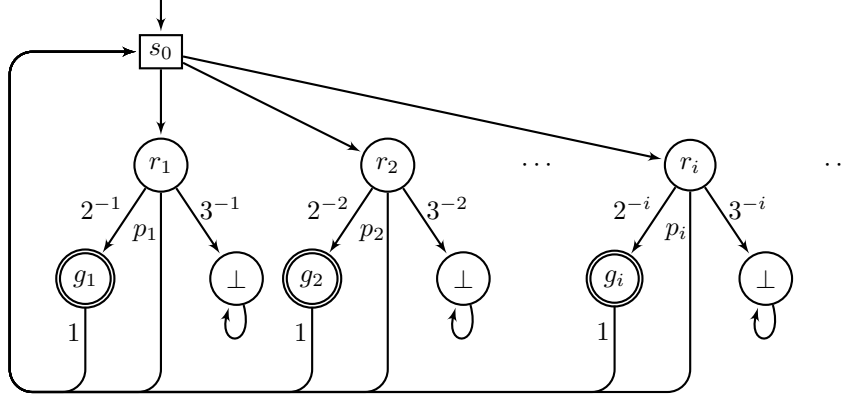


a) Finitely branching, but infinitely many controlled states.

b) Infinitely branching, but just one controlled state.

■ **Figure 1** Two MDPs where  $\varepsilon$ -optimal strategies for Büchi require infinite memory. Let  $F = \{s_0\}$  be the set of goal states. Here and throughout the paper we indicate goal states by double borders, and controlled states as rectangles.

► **Example 1.** Consider the MDPs in Figure 1. Every finite memory (FR) strategy will only attain probability 0 for Büchi in these examples [12]. However, there exists an  $\varepsilon$ -optimal Markov strategy for every  $\varepsilon > 0$ : At the  $i$ -th time that state  $s_0$  is visited, pick the successor state  $r_{i+k}$  where  $k$  is some sufficiently large number depending on  $\varepsilon$ , e.g.,  $k = \lceil \log_2(1/\varepsilon) \rceil$ . For example (b) this can easily be done with a step-counter since  $s_0$  is visited for the  $i$ -th time in step  $2(i-1)$  unless the system has reached the state  $\perp$ . For example (a), under this strategy, state  $s_0$  is visited for the  $i$ -th time in step  $\sum_{j=1}^{i-1} (k+j+1)$  unless the system has reached the state  $\perp$ . ◀



■ **Figure 2** An MDP where  $\varepsilon$ -optimal strategies for Büchi require infinite memory. The transition probability  $p_i$  stands for  $1 - 2^{-i} - 3^{-i}$ . The state  $s_0$  is the only controlled state.

► **Example 2.** Consider the MDP from Figure 2, taken from [11, Example 4.2]. Every FR-strategy attains only probability 0 of Büchi. Moreover, the strategy that, in state  $s_0$ , subsequently picks  $r_1, r_2, \dots$  also attains probability 0, unlike in Example 1. But a different infinite-memory strategy achieves a positive probability. Indeed, let  $\sigma$  be the strategy that, in  $s_0$ , picks  $r_1$   $2^1$  times and then  $r_2$   $2^2$  times and  $\dots$   $r_i$   $2^i$  times etc. This strategy  $\sigma$  achieves a positive probability of Büchi. (In more detail,  $\sigma$  achieves a positive probability of not falling in a losing sink  $\perp$ , and in almost all of the remaining runs it visits a goal state infinitely often.) Note that  $\sigma$  is a Markov strategy. ◀

**The open problem.** While the MDPs in Examples 1 and 2 require infinite memory, Markov strategies suffice for them. Such examples led to the question whether there always exists a family of  $\varepsilon$ -optimal Markov strategies for Büchi in all countably infinite MDPs.

A partial answer was given by Hill [10] (Proposition 5.1), who showed that  $\varepsilon$ -optimal Markov strategies for Büchi exist in the special case where the MDP contains only a *finite* number of controlled states. This result applies to the MDPs from Example 2 and Figure 1b), but not directly to the one in Figure 1a).

The question for general MDPs was stated as an open problem in [10] (p.158, 1.4) and mentioned again in [11] (Q1 in Section 5).

**Our contributions.** We provide a negative answer to the open problem. We construct a non-trivial example of a countable acyclic and finitely branching MDP and prove that no  $\varepsilon$ -optimal Markov strategies for Büchi exist for it (for any  $\varepsilon < 1$ ). In combination with the example from Figure 1, this shows that for general MDPs neither finite memory (FR) nor Markov strategies are sufficient.

Secondly, we provide an upper bound on the strategy complexity of Büchi. We show that for *acyclic* countable MDPs there always exist  $\varepsilon$ -optimal strategies that are deterministic and use only one bit of memory. Since every MDP can be transformed into an acyclic one by encoding a step-counter into the states, it follows that general countable MDPs have  $\varepsilon$ -optimal strategies for Büchi that are deterministic and use only a step-counter plus one extra bit of memory. Thus Markov strategies are almost, but not quite, sufficient. Table 1 summarizes these results.

$\varepsilon$ -optimal strategy for Büchi	MD	1-bit D	FR	Markov	Markov+1 bit D
Finite MDP	Y	Y	Y	Y	Y
MDP w. finitely many controlled states	N	N	N	Y	Y
Acyclic MDP	<b>N</b>	<b>Y</b>	<b>Y</b>	<b>N</b>	<b>Y</b>
General MDP	N	N	N	<b>N</b>	<b>Y</b>

■ **Table 1** Existence of various types of  $\varepsilon$ -optimal strategies for the Büchi objective, for several classes of MDPs. New results are in boldface.

## 2 Preliminaries

A *probability distribution* over a countable set  $S$  is a function  $f : S \rightarrow [0, 1]$  with  $\sum_{s \in S} f(s) = 1$ . We write  $\mathcal{D}(S)$  for the set of all probability distributions over  $S$ .

For a set  $S$  we write  $S^*$  (resp.  $S^\omega$ ) for the set of all finite (resp. infinite) sequences of elements in  $S$ . We use slightly generalized regular expressions for sets of sequences, e.g., if  $s_0 \in S$  we may write  $s_0 S^\omega$  for the set of infinite sequences starting with  $s_0$ .

**Markov decision processes.** A *Markov decision process* (MDP)  $\mathcal{M} = (S, S_\square, S_\circ, \longrightarrow, P)$  consists of a countable set  $S$  of *states*, which is partitioned into a set  $S_\square$  of *controlled states* and a set  $S_\circ$  of *random states*, a *transition relation*  $\longrightarrow \subseteq S \times S$ , and a *probability function*  $P : S_\circ \rightarrow \mathcal{D}(S)$ . We write  $s \longrightarrow s'$  if  $(s, s') \in \longrightarrow$ , and refer to  $s'$  as a *successor* of  $s$ . We assume that every state has at least one successor. The probability function  $P$  assigns to each random state  $s \in S_\circ$  a probability distribution  $P(s)$  over its (non-empty) set of successor states. A *sink* in  $\mathcal{M}$  is a subset  $T \subseteq S$  closed under the  $\longrightarrow$  relation, that is,  $s \in T$  and  $s \longrightarrow s'$  implies that  $s' \in T$ .

An MDP is *acyclic* if the underlying directed graph  $(S, \longrightarrow)$  is acyclic, i.e., there is no directed cycle. It is *finitely branching* if every state has finitely many successors and *infinitely branching* otherwise. An MDP without controlled states ( $S_\square = \emptyset$ ) is called a *Markov chain*.

**Strategies and Probability Measures.** A *run*  $\rho$  is an infinite sequence  $s_0 s_1 \dots$  of states such that  $s_i \longrightarrow s_{i+1}$  for all  $i \in \mathbb{N}$ ; write  $\rho(i) \stackrel{\text{def}}{=} s_i$  for the  $i$ -th state along  $\rho$ . A *partial run* is a finite prefix of a run. We say that (partial) run  $\rho$  *visits*  $s$  if  $s = \rho(i)$  for some  $i$ , and that  $\rho$  starts in  $s$  if  $s = \rho(0)$ .

A *strategy* is a function  $\sigma : S^* S_\square \rightarrow \mathcal{D}(S)$  that assigns to partial runs  $\rho s \in S^* S_\square$  a distribution over the successors  $\{s' \in S \mid s \longrightarrow s'\}$ . The set of all strategies in  $\mathcal{M}$  is denoted by  $\Sigma_{\mathcal{M}}$  (we omit the subscript and write  $\Sigma$  if  $\mathcal{M}$  is clear from the context). A (partial) run  $s_0 s_1 \dots$  is induced by strategy  $\sigma$  if for all  $i$  either  $s_i \in S_\square$  and  $\sigma(s_0 s_1 \dots s_i)(s_{i+1}) > 0$ , or  $s_i \in S_\circ$  and  $P(s_i)(s_{i+1}) > 0$ .

An MDP  $\mathcal{M} = (S, S_\square, S_\circ, \longrightarrow, P)$ , an initial state  $s_0 \in S$ , and a strategy  $\sigma$  induce a probability space in which the outcomes are runs starting in  $s_0$  and with measure  $\mathcal{P}_{\mathcal{M}, s_0, \sigma}$  defined as follows. It is first defined on *cylinders*  $s_0 s_1 \dots s_n S^\omega$ , where  $s_1, \dots, s_n \in S$ : if

$s_0s_1 \dots s_n$  is not a partial run induced by  $\sigma$  then  $\mathcal{P}_{\mathcal{M},s_0,\sigma}(s_0s_1 \dots s_nS^\omega) \stackrel{\text{def}}{=} 0$ . Otherwise,  $\mathcal{P}_{\mathcal{M},s_0,\sigma}(s_0s_1 \dots s_nS^\omega) \stackrel{\text{def}}{=} \prod_{i=0}^{n-1} \bar{\sigma}(s_0s_1 \dots s_i)(s_{i+1})$ , where  $\bar{\sigma}$  is the map that extends  $\sigma$  by  $\bar{\sigma}(ws) = P(s)$  for all  $ws \in S^*S_\square$ . By Carathéodory's theorem [3], this extends uniquely to a probability measure  $\mathcal{P}_{\mathcal{M},s_0,\sigma}$  on the Borel  $\sigma$ -algebra  $\mathcal{F}$  of subsets of  $s_0S^\omega$ . Elements of  $\mathcal{F}$ , i.e., measurable sets of runs, are called *events* or *objectives* here. For  $X \in \mathcal{F}$  we will write  $\bar{X} \stackrel{\text{def}}{=} s_0S^\omega \setminus X \in \mathcal{F}$  for its complement and  $\mathcal{E}_{\mathcal{M},s_0,\sigma}$  for the expectation w.r.t.  $\mathcal{P}_{\mathcal{M},s_0,\sigma}$ . We drop the indices wherever possible without introducing ambiguity.

**Strategy Classes.** Strategies are in general *randomized* (R) in the sense that they take values in  $\mathcal{D}(S)$ . A strategy  $\sigma$  is *deterministic* (D) if  $\sigma(\rho)$  is a Dirac distribution for all runs  $\rho \in S^*S_\square$ .

We formalize the amount of *memory* needed to implement strategies. Let  $\mathbf{M}$  be a countable set of memory modes, and let  $\tau : \mathbf{M} \times S \rightarrow \mathcal{D}(\mathbf{M} \times S)$  be a function that meets the following two conditions: for all modes  $\mathbf{m} \in \mathbf{M}$ ,

- for all controlled states  $s \in S_\square$ , the distribution  $\tau(\mathbf{m}, s)$  is over  $\mathbf{M} \times \{s' \mid s \rightarrow s'\}$ .
- for all random states  $s \in S_\circ$ , we have  $\sum_{\mathbf{m}' \in \mathbf{M}} \tau(\mathbf{m}, s)(\mathbf{m}', s') = P(s)(s')$ .

The function  $\tau$  together with an initial memory mode  $\mathbf{m}_0$  induce a strategy  $\sigma_\tau : S^*S_\square \rightarrow \mathcal{D}(S)$  as follows. Consider the Markov chain with the set  $\mathbf{M} \times S$  of states and the probability function  $\tau$ . A sequence  $\rho = s_0 \dots s_i$  corresponds to a set  $H(\rho) = \{(\mathbf{m}_0, s_0) \dots (\mathbf{m}_i, s_i) \mid \mathbf{m}_0, \dots, \mathbf{m}_i \in \mathbf{M}\}$  of runs in this Markov chain. Each  $\rho s \in s_0S^*S_\square$  induces a probability distribution  $\mu_{\rho s} \in \mathcal{D}(\mathbf{M})$ , the probability of being in state  $(\mathbf{m}, s)$  conditioned on having taken some partial run from  $H(\rho s)$ . We define  $\sigma_\tau$  such that  $\sigma_\tau(\rho s)(s') = \sum_{\mathbf{m}, \mathbf{m}' \in \mathbf{M}} \mu_{\rho s}(\mathbf{m}) \tau(\mathbf{m}, s)(\mathbf{m}', s')$  for all  $\rho s \in S^*S_\square$  and all  $s' \in S$ .

We say that a strategy  $\sigma$  can be *implemented* with memory  $\mathbf{M}$  if there exist  $\mathbf{m}_0 \in \mathbf{M}$  and  $\tau$  such that  $\sigma_\tau = \sigma$ . We define certain classes of strategies:

- A strategy  $\sigma$  is *finite memory* (F) if there exists a finite memory  $\mathbf{M}$  implementing  $\sigma$ .
- A strategy  $\sigma$  is *memoryless* (M) (also called *positional*) if it can be implemented with a memory of size 1. We may view M-strategies as functions  $\sigma : S_\square \rightarrow \mathcal{D}(S)$ .
- A strategy  $\sigma$  is *1-bit* if it can be implemented with a memory of size 2. Such a strategy is then determined by a function  $\tau : \{0, 1\} \times S \rightarrow \mathcal{D}(\{0, 1\} \times S)$ . Intuitively  $\tau$  uses one bit of memory to capture two different modes.
- A strategy  $\sigma$  is *Markov* if it can be implemented with the natural numbers  $\mathbb{N}$  as the memory, and a function  $\tau$  such that the distribution  $\tau(\mathbf{m}, s)$  is over  $\{\mathbf{m} + 1\} \times S$  for all  $\mathbf{m} \in \mathbf{M}$  and  $s \in S$ . Intuitively, such a strategy depends only on the the current state and the number of steps taken so far, i.e., it has access to a step-counter. We view Markov strategies as functions  $\sigma : \mathbb{N} \times S_\square \rightarrow \mathcal{D}(S)$ . Note that such a strategy is generally not finite memory.
- A strategy  $\sigma$  is *1-bit Markov* if it can be implemented with  $\mathbb{N} \times \{0, 1\}$  as the memory, and a function  $\tau$  such that the distribution  $\tau(n, b, s)$  is over  $\{n + 1\} \times \{0, 1\} \times S$  for all  $(n, b) \in \mathbf{M}$  and  $s \in S$ . We view such strategies as functions  $\sigma : \mathbb{N} \times \{0, 1\} \times S_\square \rightarrow \mathcal{D}(\{0, 1\} \times S)$ .

**Payoffs, Values, Optimality.** We are interested in strategies to maximize the expectation of a given measurable *payoff* function  $f : S^\omega \rightarrow \mathbb{R}$ , a random variable that assigns a real value to every run. The *value* of state  $s$  (w.r.t.  $f$ ) is the supremum of expected values of  $f$  over all strategies:

$$\text{val}_{\mathcal{M},f}(s) \stackrel{\text{def}}{=} \sup_{\sigma \in \Sigma} \mathcal{E}_{\mathcal{M},s,\sigma}(f),$$

For  $\varepsilon \geq 0$  and  $s \in S$ , we say that a strategy  $\sigma$  is  $\varepsilon$ -optimal iff  $\mathcal{E}_{\mathcal{M},s,\sigma}(f) \geq \text{val}_{\mathcal{M},f}(s) - \varepsilon$  and *uniformly*  $\varepsilon$ -optimal iff this holds for every  $s \in S$ . A (uniformly) 0-optimal strategy is simply called (uniformly) *optimal*.

In this paper, we will need two types of payoff functions. The first is the *total reward*, a random variable given as  $f(\rho) \stackrel{\text{def}}{=} \sum_{t=0}^{\infty} r(\rho(t))$ , where  $r : S \rightarrow \mathbb{R}$  is some given *reward* function. A useful fact [15, Theorem 7.1.9] is that if  $S$  is finite and the range of  $r$  is bounded then there exist optimal strategies (for total reward) which are memoryless and deterministic.

The second type of payoff functions we consider are those with range  $\{0, 1\}$ . Each such payoff function  $f$  uniquely identifies an objective (set of runs)  $\varphi$  by viewing  $f$  as the characteristic function of  $\varphi$ , i.e.,  $f(\rho) = 1$  if  $\rho \in \varphi$  and 0 otherwise. Then  $\mathcal{E}_{\mathcal{M},s,\sigma}(f) = \mathcal{P}_{\mathcal{M},s,\sigma}(\varphi)$ . We call this the *probability of achieving*  $\varphi$  (using strategy  $\sigma$  starting from the state  $s$ ) and simply write  $\text{val}_{\mathcal{M},\varphi}(s) = \text{val}_{\mathcal{M},f}(s) = \sup_{\sigma \in \Sigma} \mathcal{P}_{\mathcal{M},s,\sigma}(\varphi)$ .

Our main focus are *reachability* (sometimes also called *goal*) and *Büchi* objectives, which are determined by a set of states  $F \subseteq S$  and defined as follows. Let us slightly abuse notation and identify  $F$  with its characteristic function, i.e.,  $F(s) = 1$  if  $s \in F$ .

- The *reachability* objective is to visit  $F$  at least once during a run. The corresponding payoff is  $f(\rho) \stackrel{\text{def}}{=} \max_{t \in \mathbb{N}} \rho(t)$ , and we define  $\text{Goal}(F) \stackrel{\text{def}}{=} \{\rho \in S^\omega \mid \max_{t \in \mathbb{N}} F(\rho(t)) = 1\}$ ;
- The *Büchi* objective is to visit  $F$  infinitely often. The corresponding payoff function is  $f(\rho) \stackrel{\text{def}}{=} \limsup_{t \rightarrow \infty} F(\rho(t))$ , and we let  $\text{Büchi}(F) \stackrel{\text{def}}{=} \{\rho \in S^\omega \mid \limsup_{t \rightarrow \infty} F(\rho(t)) = 1\}$ .

### 3 The Lower Bound

In this section we solve Hill's problem ([10] and [11, Q1]) by exhibiting an MDP where the initial state has value 1 w.r.t. the Büchi objective, but every Markov strategy achieves this objective with probability 0. As explained in the introduction, it follows that in acyclic MDPs,  $\varepsilon$ -optimal MR-strategies are not guaranteed to exist. In fact, in the following theorem we prove the latter fact first, and subsequently generalize it to solve Hill's problem.

► **Theorem 3.** *There exists an acyclic MDP  $\mathcal{M}$ , a state  $s_0$  and a set of states  $F$  such that*

1. *for every Markov strategy  $\sigma$ , we have  $\mathcal{P}_{\mathcal{M},s_0,\sigma}(\text{Büchi}(F)) = 0$ , and*
2.  *$\text{val}_{\text{Büchi}(F)}(s_0) = 1$  and for every  $\varepsilon > 0$  there exists a deterministic 1-bit strategy  $\sigma_\varepsilon$  s.t.  $\mathcal{P}_{\mathcal{M},s_0,\sigma_\varepsilon}(\text{Büchi}(F)) \geq 1 - \varepsilon$ .*

In the remainder of this section we provide a proof sketch. The full proof is in Appendix A.

**Proof sketch for Theorem 3.** Our construction is based on an infinite MDP  $\mathcal{M}$  that consists of a chain of height- $n$  trees,  $T^n$ , for  $n \in \mathbb{N} = \{1, 2, \dots\}$ . Figure 3 depicts its initial segment  $T^1, T^2, T^3$ . Each such tree is “rooted” at a brown state on the top level, with a transition incoming from a blue state. We make use of some conventions that simplify the presentation and the analysis. In Figure 3, the different colors of the states highlight the structure of the MDP; the colors are also indicated by letters in the states: blue (L), brown (B), yellow (Y), red (R), green (G), white (W). The start state,  $s_0$ , is the blue state in the top-left corner. The controlled states are exactly the yellow states. The goal set  $F$  consists of the green states at the bottom. Two transitions emanate from each red state: a black (right) transition and a red (left) transition, both leading to the same (brown or green) state.

We consider the strengthened Büchi objective that asks to see  $F$  infinitely often and moreover that *no red transition* is taken. This corresponds exactly to the normal Büchi



objective if we redirect every red transition to an infinite (losing) chain of non-green states (not depicted in Figure 3).

We first argue that no MR-strategy achieves a positive probability of that objective. Then we show that the MDP  $\mathcal{M}$  can be modified so that no Markov strategy achieves a positive probability.

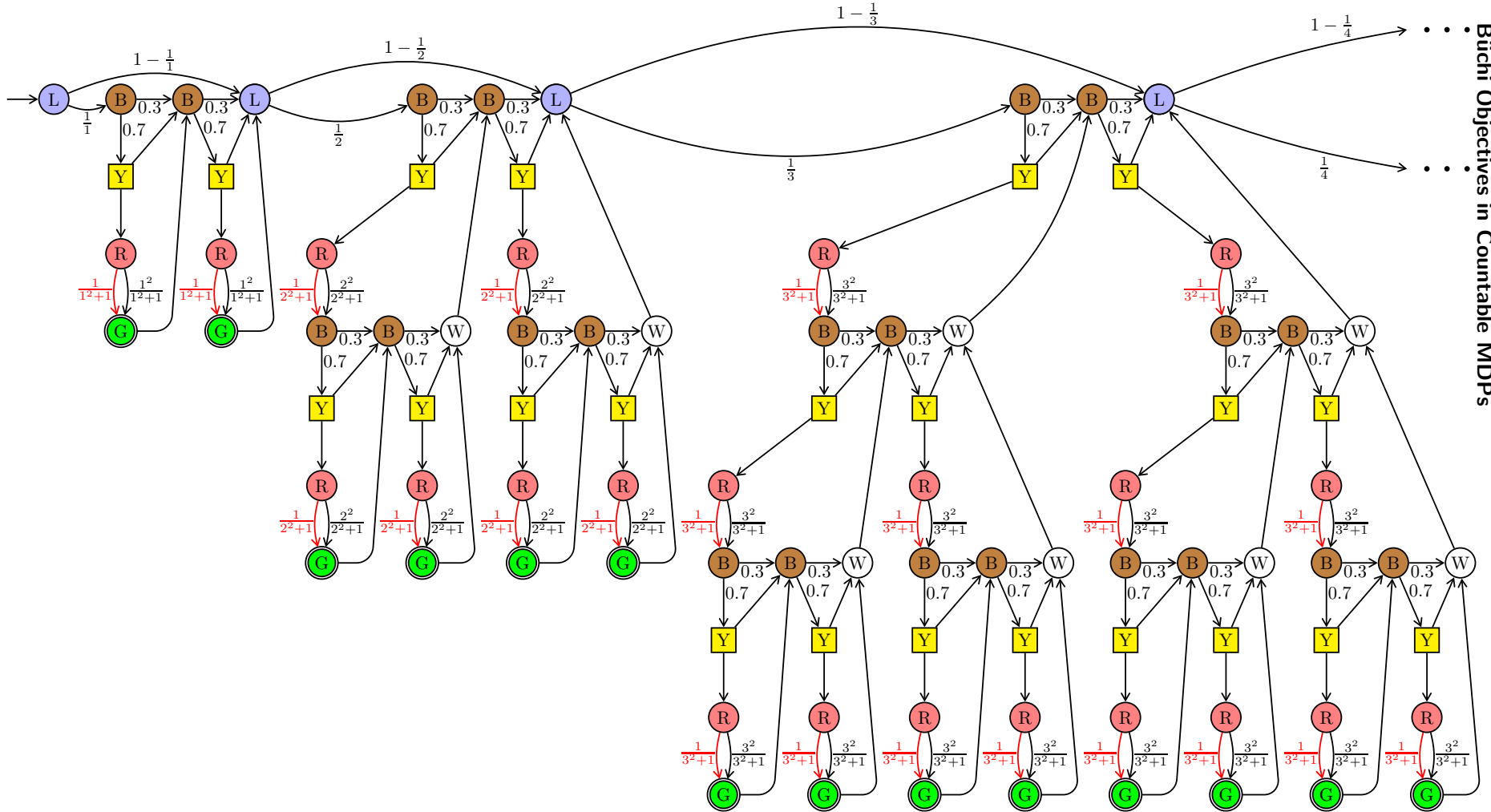
**Intuition behind the construction of  $\mathcal{M}$ .** The objective, say  $\varphi$ , of visiting infinitely many green states and no red transition creates tension between trying to visit green states and avoiding too many red states (the latter states incur a risk of taking a red transition). In the proof we need to show that no memoryless strategy strikes a good balance between these competing goals. On the one end of the spectrum, an MR-strategy might always choose the upward transition in the yellow states (which are the only controlled states). But such a strategy never visits any green state, thus clearly violates  $\varphi$ . On the other end of the spectrum lies the “greedy” MR-strategy, which always chooses the downward transition in the yellow states, in order to visit as many green states as possible. Indeed, under this strategy, let  $u_n$  denote the probability that, starting in the top-left brown state of  $T^n$ , no green state is visited in  $T^n$ . By induction (given in the appendix as part of the general proof) one can show that there is  $u < 1$  such that  $u_n \leq u$  holds for all  $n$ . Considering the probability of the transitions emanating from the blue states (at the top), the expected overall number of visited green states is at least  $\sum_{n=1}^{\infty} \frac{1}{n}(1 - u_n) \geq (1 - u) \sum_{n=1}^{\infty} \frac{1}{n} = \infty$ . It is not hard to strengthen this statement so that the greedy strategy almost surely visits infinitely many green states. So the greedy strategy satisfies one part of  $\varphi$ , but it does so at the expense of visiting many red states. Red states though are associated with a risk of taking a red transition, and it follows from the proof in the appendix that the greedy strategy almost surely ends up taking at least one (and indeed infinitely many) red transition(s).

**Good 1-bit strategies.** The two competing goals discussed in the previous paragraph can be balanced using a deterministic 1-bit strategy, which we describe in the following. This strategy,  $\sigma_1$ , sets its bit to 0 whenever a blue state (at the top) is entered. While the bit is 0, in each tree  $T^n$  it maximizes the probability of visiting a green state by choosing the downward transition in the yellow states, thus accepting a certain risk of taking a red transition. However, if and when a green state in  $T^n$  is visited, the bit is set to 1, and for the remaining sojourn in  $T^n$  the strategy  $\sigma_1$  chooses the upward transitions in the yellow states, thus avoiding any risk of a red transition in the remainder of  $T^n$ . Although  $\sigma_1$  appears to visit fewer green states than the aforementioned “greedy” MR-strategy,  $\sigma_1$  still visits infinitely many green states almost surely. This is because for each tree  $T^n$ , the two strategies have the same probability of visiting at least one green state in  $T^n$ . The strategy  $\sigma_1$  can be improved, for each  $\varepsilon > 0$ , to achieve  $\varphi$  with probability at least  $1 - \varepsilon$ , by fixing the bit to 1 in the first  $k$  trees  $T^1, \dots, T^k$ , for a  $k$  that depends on  $\varepsilon$ . Thus the first  $k$  trees are virtually skipped, eliminating the risk of taking any red transition there. In this way one can make the risk of taking a red transition arbitrarily small, while still visiting infinitely many green states with probability 1.

**No good MR-strategies.** We need to show that not only the extreme MR-strategies described above are inadequate but that every MR-strategy achieves  $\varphi$  with probability 0. To this end, for each tree  $T^n$ , define two probabilities:

- $t_n$  (for “total success”): the probability that, starting in the top-left brown state of  $T^n$ , at least one green state but no red transition is visited in  $T^n$ ;
- $d_n$  (for “death”): the probability that, starting in the top-left brown state of  $T^n$ , a red transition is visited in  $T^n$ .





■ **Figure 3** For this acyclic MDP  $\mathcal{M}$  there are  $\varepsilon$ -optimal deterministic 1-bit strategies for the Büchi objective  $\text{Büchi}(T)$  where  $T$  contains exactly all green states. No MR-strategy achieves even a positive probability.

A very technical proof shows that  $d_n \geq 0.008 \cdot t_n$  holds for all  $n$ , and this key inequality captures the inability of *any* MR-strategy to strike an adequate balance between the mentioned competing goals. Indeed, one can show that for an MR-strategy to have a positive probability of not visiting any red transition, the series  $\sum_{n=1}^{\infty} \frac{1}{n} \cdot d_n$  needs to converge; but to have a positive probability of visiting infinitely many green states, the series  $\sum_{n=1}^{\infty} \frac{1}{n} \cdot t_n$  needs to diverge (in both cases, the factor  $\frac{1}{n}$  is the probability of visiting the top-left brown node of  $T^n$ ). By the inequality above, this is impossible.

**No good Markov strategies.** For the proof of [Theorem 3](#), we also need to show that all Markov strategies achieve probability 0. To this end, we modify the MDP  $\mathcal{M}$  so that for each state, all paths from the initial state  $s_0$  to  $s$  have the same length. This can be achieved by replacing some transitions in  $\mathcal{M}$  by longer chains consisting of non-green states. This modification does not change the fact that MR-strategies achieve probability 0. But since in the new MDP each state can only be visited at a certain time, which is known a priori, a step-counter does not help. Hence all Markov strategies, like MR-strategies, achieve  $\varphi$  with probability 0. ◀

[Theorem 3](#) answers Hill's question negatively. By combining the MDP from [Theorem 3](#) with one of the MDPs from [Figure 1](#) (by adding a new initial random state that branches to the MDPs with probability  $\frac{1}{2}$  each), one can even construct a single MDP whose value w.r.t.  $\text{Büchi}(F)$  is 1, but every FR- and every Markov strategy achieves probability 0.

A slight modification of the example above yields a lower bound on the memory requirements for the almost-sure parity objective. Recall that the parity objective is defined on systems whose states are labeled by a finite set of colors  $C \stackrel{\text{def}}{=} \{1, 2, \dots, \max\} \subseteq \mathbb{N}$ , where a run is in  $\text{Parity}(C)$  iff the highest color that is seen infinitely often in the run is even.

► **Corollary 4.** *There exist an acyclic MDP  $\mathcal{M}'$  with colors  $\{1, 2, 3\}$  and a state  $s_0$  such that*

1. *for every Markov strategy  $\sigma$ , we have  $\mathcal{P}_{\mathcal{M}', s_0, \sigma}(\text{Parity}(\{1, 2, 3\})) = 0$ , and*
2. *there exists a deterministic 1-bit strategy  $\sigma'$  such that  $\mathcal{P}_{\mathcal{M}', s_0, \sigma'}(\text{Parity}(\{1, 2, 3\})) = 1$ .*

**Proof.** We obtain  $\mathcal{M}'$  by modifying the MDP  $\mathcal{M}$  from [Theorem 3](#) as follows. Label all green states in  $F$  by color 2 and the rest by color 1. Then modify each red transition to go to its target via a fresh state labeled by color 3. Clearly  $\mathcal{M}'$  is still acyclic and labeled by colors  $\{1, 2, 3\}$ .

From the proof of [Theorem 3](#) (1), under every Markov strategy in  $\mathcal{M}$  a.s. seeing infinitely many green states (in  $F$ ) implies seeing infinitely many red transitions. So in  $\mathcal{M}'$  every Markov strategy  $\sigma$  a.s. either sees color 2 only finitely often or color 3 infinitely often, thus  $\mathcal{P}_{\mathcal{M}', s_0, \sigma}(\text{Parity}(\{1, 2, 3\})) = 0$ .

From the proof of [Theorem 3](#) (2), there is a deterministic 1-bit strategy  $\sigma$  in  $\mathcal{M}$  that attains probability  $\geq 1/2$  for  $\text{Büchi}(F)$  without taking any red transition and otherwise a.s. takes a red transition. This property of  $\sigma$  holds not only when starting from  $s_0$  but from every other state as well. We obtain  $\sigma'$  in  $\mathcal{M}'$  by continuing to play  $\sigma$  even after red transitions have been taken. Under  $\sigma'$  the probability of going through infinitely many red transitions (and seeing color 3) is  $\leq (1/2)^\infty = 0$ , and the probability of seeing infinitely many states in  $F$  (with color 2) is 1. Thus  $\mathcal{P}_{\mathcal{M}', s_0, \sigma'}(\text{Parity}(\{1, 2, 3\})) = 1$ . ◀

## 4 The Upper Bound

We show that acyclic MDPs admit  $\varepsilon$ -optimal deterministic 1-bit strategies for Büchi.

We start by giving some intuition why 1 bit of memory is needed and how it is used. A step  $s' \rightarrow s''$  from some controlled state  $s'$  is *value-decreasing* iff  $\text{val}(s'') < \text{val}(s')$ . While an optimal strategy can never tolerate any value-decreasing step, an  $\varepsilon$ -optimal strategy might have to take value-decreasing steps infinitely often. The trick is to keep the collective value-loss sufficiently small ( $\leq \varepsilon$ ), while satisfying the other requirements of the objective. So the strategy needs to play ‘ever better’ (i.e., tolerate only smaller and smaller value decreases) along a run. In general this requires infinite memory, since one might re-visit the same state infinitely often and needs to choose a different transition from it every time; cf. Figure 1. However, in an acyclic MDP, with high probability, the distance to the initial state increases with the number of steps taken. Thus one can partition the state space into separate regions, depending on the distance from the initial state, and fix an acceptable rate of value-decrease for each region. Just limiting the collective value-loss is not sufficient for Büchi, one also needs to make progress and visit the set of goal states  $F$  at least once in each region. The problem is that some runs might linger in some region too long, and visit  $F$  many times, but see too many value-decreasing steps at the rate of this region. Therefore, as soon as one has visited  $F$  in some region, one should try to get to the next outer region (further away from the initial state) where the rate of value-loss is smaller. Thus one needs 1 bit of memory to record whether one has already seen  $F$  in this region. (Remember that the same state can be reached by different runs with different histories.) Just 1 bit suffices, because the probability of returning to a previous inner region (and misinterpreting the bit) can be made arbitrarily small, since the MDP is acyclic.

► **Theorem 5.** *For every acyclic countable MDP  $\mathcal{M}$ , finite set of initial states  $I$ , set of states  $F$  and  $\varepsilon > 0$ , there exists a deterministic 1-bit strategy for Büchi( $F$ ) that is  $\varepsilon$ -optimal from every  $s \in I$ .*

**Proof.** Let  $\mathcal{M} = (S, S_\square, S_\circ, \rightarrow, P)$  be an acyclic MDP,  $I \subseteq S$  a finite set of initial states and  $F \subseteq S$  a set of goal states and  $\varphi \stackrel{\text{def}}{=} \text{Büchi}(F)$  denote the Büchi objective w.r.t.  $F$ . We prove the claim for finitely branching  $\mathcal{M}$  first and transfer the result to general MDPs at the end. For every  $\varepsilon > 0$  and every  $s \in I$  there exists an  $\varepsilon$ -optimal strategy  $\sigma_s$  such that

$$\mathcal{P}_{\mathcal{M}, s, \sigma_s}(\varphi) \geq \text{val}_{\mathcal{M}, \varphi}(s) - \varepsilon. \quad (1)$$

However, the strategies  $\sigma_s$  might differ from each other and might use randomization and a large (or even infinite) amount of memory. We will construct a single deterministic strategy  $\sigma'$  that uses only 1 bit of memory such that  $\forall s \in I \mathcal{P}_{\mathcal{M}, s, \sigma'}(\varphi) \geq \text{val}_{\mathcal{M}, \varphi}(s) - 2\varepsilon$ . This proves the claim as  $\varepsilon$  can be chosen arbitrarily small.

In order to construct  $\sigma'$ , we first observe the behavior of the finitely many  $\sigma_s$  for  $s \in I$  on an infinite, increasing sequence of finite subsets of  $S$ . Based on this, we define a second stronger objective  $\varphi'$  with

$$\varphi' \subseteq \varphi, \quad (2)$$

and show that all  $\sigma_s$  attain at least  $\text{val}_{\mathcal{M}, \varphi}(s) - 2\varepsilon$  w.r.t.  $\varphi'$ , i.e.,

$$\forall s \in I \mathcal{P}_{\mathcal{M}, s, \sigma_s}(\varphi') \geq \text{val}_{\mathcal{M}, \varphi}(s) - 2\varepsilon. \quad (3)$$

We construct  $\sigma'$  as a deterministic 1-bit *optimal* strategy w.r.t.  $\varphi'$  from all  $s \in I$  and obtain

$$\begin{aligned} \mathcal{P}_{\mathcal{M}, s, \sigma'}(\varphi) &\geq \mathcal{P}_{\mathcal{M}, s, \sigma'}(\varphi') && \text{by (2)} \\ &\geq \mathcal{P}_{\mathcal{M}, s, \sigma_s}(\varphi') && \text{by optimality of } \sigma' \text{ for } \varphi' \\ &\geq \text{val}_{\mathcal{M}, \varphi}(s) - 2\varepsilon && \text{by (3).} \end{aligned}$$

**Informal outline: Behavior of  $\sigma_s$ , objective  $\varphi'$  and properties (2) and (3).** For the formal proof see Appendix B.

Let  $\text{bubble}_k(I)$  be the set of states that are reachable from some initial state in  $I$  within at most  $k$  steps. Since  $I$  is finite and  $\mathcal{M}$  is finitely branching,  $\text{bubble}_k(I)$  is finite for every  $k$ .

We define a sequence of sufficiently large and increasing numbers  $k_i$  and  $l_i$  with  $k_i < l_i < k_{i+1}$  for  $i \in \mathbb{N}$  and finite sets  $K_i \stackrel{\text{def}}{=} \text{bubble}_{k_i}(I)$  and  $L_i \stackrel{\text{def}}{=} \text{bubble}_{l_i}(I)$ . Every run from a  $s \in I$  according to  $\sigma_s$  must eventually leave each of these finite sets, because  $\mathcal{M}$  is acyclic. Moreover, we choose these numbers so that once a run has left  $L_i$  it is very unlikely to return to  $K_i$ . Let  $F_i \stackrel{\text{def}}{=} F \cap K_i \setminus L_{i-1}$ . Runs according to  $\sigma_s$  are very likely to follow a particular pattern. Let  $R_1 \stackrel{\text{def}}{=} (K_1 \setminus F_1)^* F_1$ ,  $R_2 \stackrel{\text{def}}{=} (K_2 \setminus F_2)^* F_2$  and  $R_{i+1} \stackrel{\text{def}}{=} (K_{i+1} \setminus (F_{i+1} \cup K_{i-1}))^* F_{i+1}$  for  $i \geq 2$ . We show that

$$\forall s \in I \mathcal{P}_{\mathcal{M},s,\sigma_s}(\varphi \cap \overline{R_1 R_2 \dots R_{i+1}(S \setminus K_i)^\omega}) \leq \varepsilon \quad (4)$$

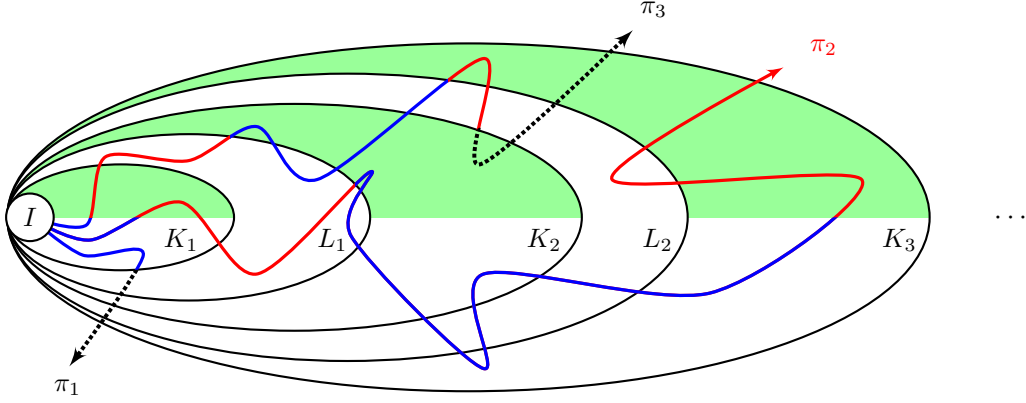
We now define the Borel objectives  $R_{\leq i} \stackrel{\text{def}}{=} R_1 R_2 \dots R_i S^\omega$  and  $\varphi' \stackrel{\text{def}}{=} \bigcap_{i \in \mathbb{N}} R_{\leq i}$ . Since  $F_i \cap F_k = \emptyset$  for  $i \neq k$  and  $\varphi'$  implies a visit to the set  $F_i$  for all  $i \in \mathbb{N}$ , we have  $\varphi' \subseteq \varphi$  and obtain (2). Using (4), we show that  $\forall s \in I \mathcal{P}_{\mathcal{M},s,\sigma_s}(\varphi') \geq \text{val}_{\mathcal{M},\varphi}(s) - 2\varepsilon$  and thus obtain (3).

**Definition of the 1-bit strategy  $\sigma'$ .** We now define a deterministic 1-bit strategy  $\sigma'$  that is optimal for objective  $\varphi'$  from every  $s \in I$ . First we define certain “suffix” objectives of  $\varphi'$ . Recall that  $R_i = (K_i \setminus (F_i \cup K_{i-2}))^* F_i$ . Let  $R_{i,j} \stackrel{\text{def}}{=} R_i R_{i+1} \dots R_j S^\omega$  and  $R_{\geq i} \stackrel{\text{def}}{=} \bigcap_{j \geq i} R_{i,j}$ . Consider the objectives  $R_{\geq i+1}$  for runs that start in states  $s' \in F_i$ . For every state  $s' \in F_i$  we consider its value w.r.t. the objective  $R_{\geq i+1}$ , i.e.,  $\text{val}_{\mathcal{M},R_{\geq i+1}}(s') \stackrel{\text{def}}{=} \sup_{\hat{\sigma}} \mathcal{P}_{\mathcal{M},s',\hat{\sigma}}(R_{\geq i+1})$ . For every  $i \geq 1$  we consider the finite subspace  $K_i \setminus K_{i-2}$ . In particular, it contains the sets  $F_{i-1}$  and  $F_i$ . We define a bounded total reward objective  $B_i$  for runs starting in  $F_{i-1}$  as follows. Runs that exit the subspace (either by leaving  $K_i$  or by visiting  $K_{i-2}$ ) before visiting  $F_i$  get reward 0. All other runs must visit  $F_i$  eventually (since  $\mathcal{M}$  is acyclic and the subspace is finite). When some run reaches the set  $F_i$  for the first time in some state  $s'$  then this run gets the reward of  $\text{val}_{\mathcal{M},R_{\geq i+1}}(s')$ . Using [15, Theorem 7.1.9], we show that there exists a uniform optimal MD-strategy  $\sigma_i$  for  $B_i$  on  $K_i \setminus K_{i-2}$  in  $\mathcal{M}$ .

We now define  $\sigma'$  by combining different MD-strategies  $\sigma_i$ , depending on the current state and on the value of the 1-bit memory. The intuition is that the strategy  $\sigma'$  has two modes: normal-mode and next-mode. In a state  $s' \in K_i \setminus K_{i-1}$ , if the memory is  $i \pmod{2}$  then the strategy is in normal-mode and plays towards reaching  $F_i$ . Otherwise, the strategy is in next-mode and plays towards reaching  $F_{i+1}$ .

Initially  $\sigma'$  starts in a state  $s \in I$  with the 1-bit memory set to 1. We define the behavior of  $\sigma'$  in a state  $s' \in K_i \setminus K_{i-1}$  for every  $i \geq 1$ . If the 1-bit memory is  $i \pmod{2}$  and  $s' \notin F_i$  then  $\sigma'$  plays like  $\sigma_i$ . (Intuitively, one plays towards  $F_i$ , since one has not yet visited it.) If the 1-bit memory is  $i \pmod{2}$  and  $s' \in F_i$  then the 1-bit memory is set to  $(i+1) \pmod{2}$ , and  $\sigma'$  plays like  $\sigma_{i+1}$ . (Intuitively, one records the fact that one has already seen  $F_i$  and then targets the next set  $F_{i+1}$ .) If the 1-bit memory is  $(i+1) \pmod{2}$  then  $\sigma'$  plays like  $\sigma_{i+1}$ . (Intuitively, one plays towards  $F_{i+1}$ , since one has already visited  $F_i$ .)

Observe that if a run according to  $\sigma'$  exits some set  $K_i$  (and thus enters  $K_{i+1} \setminus K_i$ ) with the bit still set to  $i \pmod{2}$  (normal-mode) then this run has not visited  $F_i$  and thus does not satisfy the objective  $\varphi'$ . (Or the same has happened earlier for some  $j < i$ , in which case also the objective  $\varphi'$  is violated.) An example is the run  $\pi_1$  in Figure 4. However, if a run according to  $\sigma'$  exits some set  $K_i$  (and thus enters  $K_{i+1} \setminus K_i$ ) with the bit set to  $(i+1) \pmod{2}$  (thus  $\sigma_{i+1}$  in next-mode) then in the new set  $K_{i'} \setminus K_{i'-1}$  with  $i' = i+1$  the bit is set to  $i' \pmod{2}$  and  $\sigma'$  continues to play like  $\sigma_{i+1}$  in normal-mode. Even if this run returns (temporarily) to  $K_i$  (but not to  $K_{i-1}$ ) the strategy  $\sigma'$  continues to play like  $\sigma_{i+1}$  in



■ **Figure 4** Memory updates along runs  $\pi_1, \pi_2, \pi_3$ , drawn in blue while the memory-bit is one and in red while the bit is zero. The green region in  $K_1$  is  $F_1$ , and for all  $i \geq 2$ , the green region in  $K_i \setminus L_{i-1}$  is  $F_i$ . Both  $\pi_1$  and  $\pi_3$  violate  $\varphi'$  and are drawn as dotted lines once they do.

next-mode. An example is the run  $\pi_2$  in Figure 4. Finally, if a run returns to  $K_{i-1}$  after having visited  $F_i$  then it fails the objective  $\varphi'$ , e.g., run  $\pi_3$  in Figure 4.

**The 1-bit strategy  $\sigma'$  is optimal for  $\varphi'$  from every  $s \in I$ .** Let  $s \in I$  be arbitrary. For a given run from  $s$ , let  $\text{firstin}(F_i)$  be the first state  $s'$  in  $F_i$  that is visited (if any). We define a bounded reward objective  $B'_i$  for runs starting at  $s$  as follows. Every run that does not satisfy the objective  $R_{\leq i}$  gets assigned reward 0. Otherwise, consider a run from  $s$  that satisfies  $R_{\leq i}$ . When this run reaches the set  $F_i$  for the first time in some state  $s'$  then this run gets a reward of  $\text{val}_{\mathcal{M}, R_{\geq i+1}}(s')$ . Note that this reward is  $\leq 1$ .

We show that for all  $i \in \mathbb{N}$

$$\text{val}_{\mathcal{M}, \varphi'}(s) = \text{val}_{\mathcal{M}, B'_i}(s) \quad (5)$$

Towards the  $\geq$  inequality, let  $\hat{\sigma}$  be an  $\hat{\epsilon}$ -optimal strategy for  $B'_i$  from  $s$ . We define the strategy  $\hat{\sigma}'$  to play like  $\hat{\sigma}$  until a state  $s' \in F_i$  is reached and then to switch to some  $\hat{\epsilon}$ -optimal strategy for objective  $R_{\geq i+1}$  from  $s'$ . Every run from  $s$  that satisfies  $\varphi'$  can be split into parts, before and after the first visit to the set  $F_i$ , i.e.,  $\varphi' = \{w_1 s' w_2 \mid w_1 s' \in R_{\leq i}, s' \in F_i, s' w_2 \in R_{\geq i+1}\}$ . Therefore we obtain that  $\mathcal{P}_{\mathcal{M}, s, \hat{\sigma}'}(\varphi') \geq \mathcal{E}_{\mathcal{M}, s, \hat{\sigma}}(B'_i) - \hat{\epsilon} \geq \text{val}_{\mathcal{M}, B'_i}(s) - 2\hat{\epsilon}$ . Since this holds for every  $\hat{\epsilon} > 0$ , we obtain  $\text{val}_{\mathcal{M}, \varphi'}(s) \geq \text{val}_{\mathcal{M}, B'_i}(s)$ .

Towards the  $\leq$  inequality, let  $\hat{\sigma}$  be any strategy for  $\varphi'$  from  $s$ . We have  $\mathcal{P}_{\mathcal{M}, s, \hat{\sigma}}(\varphi') \leq \sum_{s' \in F_i} \mathcal{P}_{\mathcal{M}, s, \hat{\sigma}}(R_{\leq i} \cap \text{firstin}(F_i) = s') \cdot \text{val}_{\mathcal{M}, R_{\geq i+1}}(s') = \mathcal{E}_{\mathcal{M}, s, \hat{\sigma}}(B'_i)$ . Thus  $\text{val}_{\mathcal{M}, \varphi'}(s) \leq \text{val}_{\mathcal{M}, B'_i}(s)$ . Together we obtain (5).

For all  $i \in \mathbb{N}$  and every state  $s' \in F_i$  we show that

$$\text{val}_{\mathcal{M}, R_{\geq i+1}}(s') = \text{val}_{\mathcal{M}, B_{i+1}}(s') \quad (6)$$

Towards the  $\geq$  inequality, let  $\hat{\sigma}$  be an  $\hat{\epsilon}$ -optimal strategy for  $B_{i+1}$  from  $s' \in F_i$ . We define the strategy  $\hat{\sigma}'$  to play like  $\hat{\sigma}$  until a state  $s'' \in F_{i+1}$  is reached and then to switch to some  $\hat{\epsilon}$ -optimal strategy for objective  $R_{\geq i+2}$  from  $s''$ . We have that  $\mathcal{P}_{\mathcal{M}, s', \hat{\sigma}'}(R_{\geq i+1}) \geq \mathcal{E}_{\mathcal{M}, s', \hat{\sigma}}(B_{i+1}) - \hat{\epsilon} \geq \text{val}_{\mathcal{M}, B_{i+1}}(s) - 2\hat{\epsilon}$ . Since this holds for every  $\hat{\epsilon} > 0$ , we obtain  $\text{val}_{\mathcal{M}, R_{\geq i+1}}(s') \geq \text{val}_{\mathcal{M}, B_{i+1}}(s')$ .

Towards the  $\leq$  inequality, let  $\hat{\sigma}$  be any strategy for  $R_{\geq i+1}$  from  $s' \in F_i$ . We have

$$\begin{aligned} \mathcal{P}_{\mathcal{M}, s', \hat{\sigma}}(R_{\geq i+1}) &\leq \sum_{s'' \in F_{i+1}} \mathcal{P}_{\mathcal{M}, s', \hat{\sigma}}(R_{i+1} S^\omega \cap \text{firstin}(F_{i+1}) = s'') \cdot \text{val}_{\mathcal{M}, R_{\geq i+2}}(s'') \\ &= \mathcal{E}_{\mathcal{M}, s', \hat{\sigma}}(B_{i+1}). \end{aligned}$$

Thus  $\text{val}_{\mathcal{M}, R_{\geq i+1}}(s') \leq \text{val}_{\mathcal{M}, B_{i+1}}(s')$ . Together we obtain (6).

We show, by induction on  $i$ , that  $\sigma'$  is optimal for  $B'_i$  for all  $i \in \mathbb{N}$  from start state  $s$ , i.e.,

$$\mathcal{E}_{\mathcal{M}, s, \sigma'}(B'_i) = \text{val}_{\mathcal{M}, B'_i}(s) \quad (7)$$

In the base case of  $i = 1$  we have that  $B'_1 = B_1$ . The strategy  $\sigma'$  plays  $\sigma_1$  until reaching  $F_1$ , which is optimal for objective  $B_1$  and thus optimal for  $B'_1$ . For the induction step we assume (IH) that  $\sigma'$  is optimal for  $B'_i$ .

$$\begin{aligned} \text{val}_{\mathcal{M}, B'_{i+1}}(s) &= \text{val}_{\mathcal{M}, B'_i}(s) && \text{by (5)} \\ &= \mathcal{E}_{\mathcal{M}, s, \sigma'}(B'_i) && \text{by (IH)} \\ &= \sum_{s' \in F_i} \mathcal{P}_{\mathcal{M}, s, \sigma'}(R_{\leq i} \cap \text{firstin}(F_i) = s') \cdot \text{val}_{\mathcal{M}, R_{\geq i+1}}(s') && \text{by def. of } B'_i \\ &= \sum_{s' \in F_i} \mathcal{P}_{\mathcal{M}, s, \sigma'}(R_{\leq i} \cap \text{firstin}(F_i) = s') \cdot \text{val}_{\mathcal{M}, B_{i+1}}(s') && \text{by (6)} \\ &= \sum_{s' \in F_i} \mathcal{P}_{\mathcal{M}, s, \sigma'}(R_{\leq i} \cap \text{firstin}(F_i) = s') \cdot \mathcal{E}_{\mathcal{M}, s', \sigma_{i+1}}(B_{i+1}) && \text{opt. of } \sigma_{i+1} \text{ for } B_{i+1} \\ &= \mathcal{E}_{\mathcal{M}, s, \sigma'}(B'_{i+1}) && \text{by def. of } \sigma' \text{ and } B'_{i+1} \end{aligned}$$

So  $\sigma'$  attains the value  $\text{val}_{\mathcal{M}, B'_{i+1}}(s)$  of the objective  $B'_{i+1}$  from  $s$  and is optimal. Thus (7).

Now we show that  $\sigma'$  performs well on the objectives  $R_{\leq i}$  for all  $i \in \mathbb{N}$ .

$$\mathcal{P}_{\mathcal{M}, s, \sigma'}(R_{\leq i}) \geq \text{val}_{\mathcal{M}, \varphi'}(s) \quad (8)$$

We have

$$\begin{aligned} \mathcal{P}_{\mathcal{M}, s, \sigma'}(R_{\leq i}) &\geq \mathcal{E}_{\mathcal{M}, s, \sigma'}(B'_i) \quad \text{since } B'_i \text{ gives rewards 0 for runs } \notin R_{\leq i} \text{ and } \leq 1 \text{ otherwise} \\ &= \text{val}_{\mathcal{M}, B'_i}(s) \quad \text{by (7)} \\ &= \text{val}_{\mathcal{M}, \varphi'}(s) \quad \text{by (5)} \end{aligned}$$

So we get (8). Now we are ready to prove the optimality of  $\sigma'$  for  $\varphi'$  from  $s$ .

$$\begin{aligned} \mathcal{P}_{\mathcal{M}, s, \sigma'}(\varphi') &= \mathcal{P}_{\mathcal{M}, s, \sigma'}(\cap_{i \in \mathbb{N}} R_{\leq i}) && \text{by def. of } \varphi' \\ &= \lim_{i \rightarrow \infty} \mathcal{P}_{\mathcal{M}, s, \sigma'}(R_{\leq i}) && \text{by continuity of measures from above} \\ &\geq \lim_{i \rightarrow \infty} \text{val}_{\mathcal{M}, \varphi'}(s) && \text{by (8)} \\ &= \text{val}_{\mathcal{M}, \varphi'}(s) \end{aligned}$$

**From finitely to infinitely branching MDPs.** Encode infinite branching into finite branching like in Figure 1, apply the above result to obtain a 1-bit strategy for the finitely branching version, and then transform this strategy back into a 1-bit strategy for the original MDP.  $\blacktriangleleft$

Now we show our upper bound on the strategy complexity of Büchi for general MDPs.

► **Theorem 6.** *For every countable MDP  $\mathcal{M}$ , finite set of initial states  $I$ , set of states  $F$  and  $\varepsilon > 0$ , there exists a deterministic 1-bit Markov strategy for  $\text{Büchi}(F)$  that is  $\varepsilon$ -optimal from every  $s \in I$ .*

**Proof.** Encode a step-counter into the states to obtain an acyclic MDP, apply Theorem 5 to obtain an  $\varepsilon$ -optimal deterministic 1-bit strategy for it, and then transform this strategy back into an  $\varepsilon$ -optimal deterministic 1-bit Markov strategy in the original MDP. ◀

---

## References

---

- 1 Pieter Abbeel and Andrew Y. Ng. Learning first-order Markov models for control. In *Advances in Neural Information Processing Systems 17*, pages 1–8. MIT Press, 2004.
- 2 Christel Baier and Joost-Pieter Katoen. *Principles of Model Checking*. MIT Press, 2008.
- 3 P. Billingsley. *Probability and Measure*. Wiley, 1995. Third Edition.
- 4 Vincent D. Blondel and John N. Tsitsiklis. A survey of computational complexity results in systems and control. *Automatica*, 36(9):1249–1274, 2000.
- 5 Nicole Bäuerle and Ulrich Rieder. *Markov Decision Processes with Applications to Finance*. Springer-Verlag Berlin Heidelberg, 2011.
- 6 K. Chatterjee, L. de Alfaro, and T. Henzinger. Trading memory for randomness. In *Annual Conference on Quantitative Evaluation of Systems*, pages 206–217. IEEE Computer Society Press, 2004.
- 7 K. Chatterjee and T. Henzinger. A survey of stochastic  $\omega$ -regular games. *Journal of Computer and System Sciences*, 78(2):394–413, 2012.
- 8 K. Chatterjee, M. Jurdziński, and T. Henzinger. Quantitative stochastic parity games. In *Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 121–130. Society for Industrial and Applied Mathematics, 2004.
- 9 Edmund M. Clarke, Thomas A. Henzinger, Helmut Veith, and Roderick Bloem, editors. *Handbook of Model Checking*. Springer, 2018.
- 10 Theodore Preston Hill. On the existence of good Markov strategies. *Transactions of the American Mathematical Society*, 247:157–176, 1979.
- 11 Theodore Preston Hill. Goal problems in gambling theory. *Revista de Matemática: Teoría y Aplicaciones*, 6(2):125–132, 1999.
- 12 Stefan Kiefer, Richard Mayr, Mahsa Shirmohammadi, and Dominik Wojtczak. Parity Objectives in Countable MDPs. In *Annual IEEE Symposium on Logic in Computer Science*, 2017.
- 13 J. Krčál. Determinacy and Optimal Strategies in Stochastic Games. Master’s thesis, Masaryk University, School of Informatics, 2009.
- 14 D. Ornstein. On the existence of stationary optimal strategies. *Proceedings of the American Mathematical Society*, 20:563–569, 1969.
- 15 Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., 1st edition, 1994.
- 16 Manfred Schäl. Markov decision processes in finance and dynamic options. In *Handbook of Markov Decision Processes*, pages 461–487. Springer, 2002.
- 17 Olivier Sigaud and Olivier Buffet. *Markov Decision Processes in Artificial Intelligence*. John Wiley & Sons, 2013.
- 18 R.S. Sutton and A.G Barto. *Reinforcement Learning: An Introduction*. Adaptive Computation and Machine Learning. MIT Press, 2018.



## A The Lower Bound: Full Details

► **Theorem 3.** *There exists an acyclic MDP  $\mathcal{M}$ , a state  $s_0$  and a set of states  $F$  such that*

1. *for every Markov strategy  $\sigma$ , we have  $\mathcal{P}_{\mathcal{M}, s_0, \sigma}(\text{Büchi}(F)) = 0$ , and*
2.  *$\text{val}_{\text{Büchi}(F)}(s_0) = 1$  and for every  $\varepsilon > 0$  there exists a deterministic 1-bit strategy  $\sigma_\varepsilon$  s.t.  $\mathcal{P}_{\mathcal{M}, s_0, \sigma_\varepsilon}(\text{Büchi}(F)) \geq 1 - \varepsilon$ .*

We follow the proof sketch from the main body and first argue that, in the MDP  $\mathcal{M}$  from Figure 3, no MR-strategy achieves a positive probability for the objective of visiting  $F$  infinitely often and taking no red transition. Indeed, given an MR-strategy and a tree, we define two probabilities:

- $s$  (for “survival”): the probability that, starting in the top-left brown state, no red transition in the tree is visited;
- $t$  (for “total success”): the probability that, starting in the top-left brown state, at least one green state but no red transition in the tree is visited.

Trivially,  $t \leq s$ . A key lemma is the following.

► **Lemma 7.** *Write  $p \stackrel{\text{def}}{=} 0.7$ . For every MR-strategy  $\sigma$  and every  $n \in \mathbb{N}$ , the tree  $T^n$  satisfies:*

$$s \leq a^{qtn^2},$$

where  $a = 1 - \frac{1}{n^2+1}$  and  $q = \frac{1}{9}(1-p)$ .

**Proof.** Fix any MR-strategy  $\sigma$  and any  $n \in \mathbb{N}$ . For each  $k \in \{0, \dots, n\}$ , the tree  $T^n$  has  $2^{n-k}$  height- $k$  subtrees, for which we can define  $s, t$  analogously. We claim: for all  $k \in \{0, \dots, n\}$  the probabilities  $s, t$  in every height- $k$  subtree of  $T^n$  satisfy

$$s \leq a^{(qt + \frac{1}{2}qt^2)k^2}, \tag{9}$$

where  $a = 1 - \frac{1}{n^2+1}$  and  $q = \frac{1}{9}(1-p)$ . Note that the claim (for  $k = n$ ) implies the lemma.

We prove the claim by induction on  $k$ . For the base case,  $k = 0$ , note that each height-0 subtree of  $T^n$  consists of only a single green state. Hence  $s = t = 1$ , so the claim holds for  $k = 0$ . For the inductive step, let  $k \in \{1, \dots, n\}$  and consider a height- $k$  subtree, say  $T$ , of  $T^n$ . Let  $T_0, T_1$  be the left and the right subtree of  $T$ , respectively; they have height  $k-1$ . In the two (yellow) topmost controlled states in  $T$ , the MR-strategy  $\sigma$  chooses probabilities to visit  $T_0, T_1$ , respectively. Taking into account the two brown random states at the top, the probabilities to visit  $T_0, T_1$  are  $p_0, p_1 \leq p$ , respectively. In  $T_0, T_1$ , the strategy  $\sigma$  employs MR-strategies that achieve probabilities  $s_0, t_0$  and  $s_1, t_1$ , respectively, where  $s_i, t_i$  are defined in the obvious way for  $T_i$ . By the induction hypothesis we have

$$s_i \leq a^{qt_i(1 + \frac{1}{2}t_i)(k-1)^2} \quad \text{for } i \in \{0, 1\}. \tag{10}$$

By the structure of the MDP  $\mathcal{M}$  we have:

$$s = (1 - p_0 + p_0as_0)(1 - p_1 + p_1as_1) \tag{11}$$

$$\begin{aligned} t &= p_0at_0(1 - p_1 + p_1as_1) + p_1at_1(1 - p_0 + p_0as_0) - p_0at_0p_1at_1 \\ &\leq p_0at_0 + p_1at_1 - p_0at_0p_1at_1 \end{aligned} \tag{12}$$

By combining (10) and (11) we obtain:

$$s \leq \prod_{i=0}^1 \left( 1 - p_i + p_i a^{1 + qt_i(1 + \frac{1}{2}t_i)(k^2 - 2k)} \right) \tag{13}$$

On the other hand, from (12) we obtain:

$$\begin{aligned} qt + \frac{1}{2}qt^2 &\leq q(p_0t_0 + p_1t_1 - p_0at_0p_1at_1) + \frac{1}{2}q(p_0at_0 + p_1at_1)^2 \\ &\leq p_0qt_0 \left(1 + \frac{1}{2}p_0t_0\right) + p_1qt_1 \left(1 + \frac{1}{2}p_1t_1\right) \end{aligned} \quad (14)$$

Further we have:

$$\ln \frac{1}{a} \leq \frac{1}{a} - 1 = \frac{n^2 + 1}{n^2} - 1 = \frac{1}{n^2} \leq \frac{1}{k^2} \quad (15)$$

Let  $i \in \{0, 1\}$ . By (15) we have:

$$\left(\ln \frac{1}{a}\right) p_i q t_i \left(1 + \frac{1}{2}p_i t_i\right) k^2 \leq q \left(1 + \frac{1}{2}\right) \leq \frac{1}{9} \cdot \frac{3}{2} < \frac{1}{2}$$

Hence, using a bound on the exponential function (Lemma 8 below), we obtain:

$$\begin{aligned} a^{p_i q t_i (1 + \frac{1}{2}p_i t_i) k^2} &= e^{-p_i (\ln \frac{1}{a}) q t_i (1 + \frac{1}{2}p_i t_i) k^2} \\ &\geq 1 - p_i + p_i e^{-(\ln \frac{1}{a}) q t_i (1 + \frac{1}{2}p_i t_i) k^2 - (\ln \frac{1}{a})^2 q^2 t_i^2 (1 + \frac{1}{2}p_i t_i)^2 k^4} \\ &\stackrel{(15)}{\geq} 1 - p_i + p_i e^{-(\ln \frac{1}{a}) q t_i (1 + \frac{1}{2}p_i t_i) k^2 - (\ln \frac{1}{a}) \frac{9}{4} q^2 t_i^2 k^2} \\ &= 1 - p_i + p_i a^{q t_i k^2 + q(\frac{1}{2}p_i + \frac{9}{4}q) t_i^2 k^2} \end{aligned}$$

By combining this inequality with (14) we obtain:

$$a^{(qt + \frac{1}{2}qt^2)k^2} \geq \prod_{i=0}^1 \left(1 - p_i + p_i a^{q t_i k^2 + q(\frac{1}{2}p_i + \frac{9}{4}q) t_i^2 k^2}\right)$$

Considering (13), we see that, in order to prove (9), it suffices to prove

$$1 + q t_i \left(1 + \frac{1}{2}t_i\right) (k^2 - 2k) \geq q t_i k^2 + q \left(\frac{1}{2}p_i + \frac{9}{4}q\right) t_i^2 k^2 \quad \text{for } i \in \{0, 1\}.$$

This inequality is equivalent to:

$$\begin{aligned} 1 + q t_i k \left(\left(\frac{1}{2} - \frac{1}{2}p_i - \frac{9}{4}q\right) t_i k - 2\left(1 + \frac{1}{2}t_i\right)\right) &\geq 0 \\ \iff 1 + q t_i k \left(\left(\frac{1}{2}(1-p) - \frac{9}{4}q\right) t_i k - 3\right) &\geq 0 \\ \iff 1 + \frac{1}{9}(1-p) t_i k \left(\frac{1}{4}(1-p) t_i k - 3\right) &\geq 0 \\ \iff \left(\frac{1}{6}(1-p) t_i k - 1\right)^2 &\geq 0 \end{aligned}$$

The left-hand side is a square, hence nonnegative. This completes the induction proof. ◀

The following elementary lemma from calculus was used in the proof of Lemma 7.

► **Lemma 8.** For every  $r \geq 0$  and  $x \in [0, \frac{1}{2}]$  we have  $e^{-rx} \geq 1 - r + re^{-x-x^2}$ .

**Proof.** Let  $r \geq 0$  and  $x \in [0, \frac{1}{2}]$ . As  $1 + y \leq e^y$  holds for all  $y$ , we have:

$$1 - r + re^{-x-x^2} = 1 - r \left(1 - e^{-x-x^2}\right) \leq e^{-r(1-e^{-x-x^2})}$$

Hence it suffices to prove that  $x \leq 1 - e^{-x-x^2}$ , which is equivalent to  $\ln(1-x) + x + x^2 \geq 0$ . To prove the latter inequality, define  $f(y) \stackrel{\text{def}}{=} \ln(1-y) + y + y^2$ . Then we have  $f(0) = 0$  and

$$f'(y) = -\frac{1}{1-y} + 1 + 2y = \frac{-1 + 1 - y + 2y - 2y^2}{1-y} = \frac{y(1-2y)}{1-y} \geq 0 \text{ for } y \in \left[0, \frac{1}{2}\right].$$

By the fundamental theorem of calculus, it follows  $f(x) = f(0) + \int_0^x f'(y) dy \geq 0$ .  $\blacktriangleleft$

► **Lemma 9.** Consider the acyclic MDP  $\mathcal{M}$  shown in Figure 3. Let  $\varphi$  be the objective of visiting infinitely many green states and no red transition.

1. For every MR-strategy  $\sigma$ , we have  $\mathcal{P}_{\mathcal{M}, s_0, \sigma}(\varphi) = 0$ .
2.  $\text{val}_{\varphi}(s_0) = 1$  and for every  $\varepsilon > 0$  there exists a deterministic 1-bit strategy  $\sigma_\varepsilon$  s.t.  $\mathcal{P}_{\mathcal{M}, s_0, \sigma_\varepsilon}(\varphi) \geq 1 - \varepsilon$ .

**Proof.** First we prove item 1. Fix any MR-strategy  $\sigma$ . For each  $n \in \mathbb{N}$ , let  $s_n, t_n$  denote the probabilities  $s, t$  for the tree  $T^n$  under  $\sigma$ . Define also  $d_n \stackrel{\text{def}}{=} 1 - s_n$  (for “death”), which is the probability of taking at least one red transition starting in the top-left brown state of  $T^n$ . For the following estimate, observe that we have

$$\left(1 - \frac{1}{x+1}\right)^x = e^{x \ln(1 - \frac{1}{x+1})} \leq e^{-\frac{x}{x+1}} \leq e^{-\frac{1}{2}} \quad \text{for } x \geq 1. \quad (16)$$

By Lemma 7 we have for every  $n$ :

$$d_n = 1 - s_n \geq 1 - \left(1 - \frac{1}{n^2 + 1}\right)^{qt_n n^2} \stackrel{(16)}{\geq} 1 - e^{-\frac{1}{2}qt_n} \geq \frac{1}{4}qt_n, \quad (17)$$

where the last inequality follows from the fact that  $e^{-x} \leq 1 - \frac{1}{2}x$  holds for  $x \in [0, 1]$ .

Denote by  $G_n$  the indicator random variable such that

- $G_n = 1$  if the top-left brown state of  $T^n$  is visited (coming from the previous blue state) and at least one green state in  $T^n$  but no red transition in  $T^n$  is visited;
- $G_n = 0$  otherwise.

Considering that the probability of visiting the top-left brown state of  $T^n$  is  $\frac{1}{n}$ , we have  $\mathcal{E}G_n = \frac{1}{n} \cdot t_n$ , where  $\mathcal{E}$  denotes expectation.

If  $\sigma$  visits at least one red transition in  $\mathcal{M}$  almost surely then the probability of  $\varphi$  is 0. Therefore, suppose  $\sigma$  achieves a positive probability,  $\bar{r} > 0$ , of visiting no red transition. Since  $0 < \bar{r} = \prod_{n=1}^{\infty} (1 - \frac{1}{n} \cdot d_n)$ , the series  $\sum_{n=1}^{\infty} \frac{1}{n} \cdot d_n$  converges. Thus:

$$\mathcal{E} \sum_{n=1}^{\infty} G_n = \sum_{n=1}^{\infty} \mathcal{E}G_n = \sum_{n=1}^{\infty} \frac{1}{n} \cdot t_n \stackrel{(17)}{\leq} \frac{4}{q} \cdot \sum_{n=1}^{\infty} \frac{1}{n} \cdot d_n < \infty$$

It follows that the probability that  $\sum_{n=1}^{\infty} G_n$  diverges is 0. But on  $\varphi$  the series  $\sum_{n=1}^{\infty} G_n$  diverges. Hence the probability of  $\varphi$  is 0. This completes the proof of item 1.

Towards item 2, we first define a suitable strategy,  $\sigma$ , that achieves a positive value (i.e.,  $\mathcal{P}_{\mathcal{M}, s_0, \sigma}(\varphi) > 0$ ) and then improve it to obtain  $\varepsilon$ -optimal strategies  $\sigma_\varepsilon$ .

The strategy  $\sigma$  acts independently in each tree  $T^n$ . In each tree  $T^n$  the strategy  $\sigma$  maximizes the probability of visiting exactly one green state. To this end, as long as  $\sigma$  has not yet visited a green state in  $T^n$ , it chooses the downward transition emanating from the yellow controlled states; as soon as a green state in  $T^n$  has been visited,  $\sigma$  chooses the upward transition emanating from the yellow controlled states, thus avoiding any further visit of a green state or a red transition in  $T^n$ . This is a 1-bit strategy, as  $\sigma$  remembers only whether a green state has already been visited in the current tree  $T^n$ . The bit is reset whenever a new tree is entered.

Next we show that  $\sigma$  visits infinitely many green states with probability 1. Let  $u_n$  denote the probability that, starting in the top-left brown state of  $T^n$ , no green state is visited in  $T^n$ . Define  $u_0 = 0$ . Since red or non-red transitions are unimportant for the current considerations, all height- $n$  trees in  $\mathcal{M}$  have the same structure, even when they are subtrees of different  $T^m$ . Therefore we have:

$$u_n = (pu_{n-1} + 1 - p)^2$$

Since the function  $f(x) \stackrel{\text{def}}{=} (px + 1 - p)^2$  is monotone on  $[0, 1]$ , the sequence  $(u_n)_n$  is non-decreasing and thus converges to the smaller fixed point,  $u$ , of  $f$ . Hence,

$$0 \leq u_n \leq u = f(u) = \left(\frac{1-p}{p}\right)^2 < 1 \quad \text{for all } n \in \mathbb{N} \cup \{0\}. \quad (18)$$

It follows that we have

$$\sum_{n=k}^{\infty} \frac{1}{n}(1 - u_n) \geq (1 - u) \sum_{n=k}^{\infty} \frac{1}{n} = \infty \quad \text{for all } k \in \mathbb{N}. \quad (19)$$

For every  $k \in \mathbb{N}$ , the probability that, starting in the blue state directly before  $T^k$ , no green state in  $T^k, T^{k+1}, \dots$  is visited is

$$\prod_{n=k}^{\infty} \left( \frac{1}{n} \cdot u_n + \left(1 - \frac{1}{n}\right) \right) = \prod_{n=k}^{\infty} \left( 1 - \frac{1}{n}(1 - u_n) \right) \stackrel{\text{by (19)}}{=} 0.$$

It follows that  $\sigma$  visits infinitely many green states with probability 1.

It now suffices to show that, with positive probability,  $\sigma$  visits no red transition. Let  $v_n$  denote the expectation, starting in the top-left brown state of  $T^n$ , of the number of red *states* (not red *transitions*) that are visited in  $T^n$ . Define  $v_0 \stackrel{\text{def}}{=} 0$ . Since red or non-red transitions are unimportant for the current considerations, all height- $n$  trees in  $\mathcal{M}$  have the same structure, even when they are subtrees of different  $T^m$ . Therefore we have:

$$v_n = p(1 + v_{n-1} + u_{n-1}p(1 + v_{n-1})) + (1 - p)p(1 + v_{n-1}) \quad (20)$$

We prove by induction that  $v_n \leq n$  holds for all  $n \in \mathbb{N} \cup \{0\}$ . The base case,  $n = 0$ , holds by the definition of  $v_0$ . For the inductive step, let  $n \geq 1$ . We have:

$$\begin{aligned} v_n &\leq p(n + u_{n-1}pn) + (1 - p)pn && \text{by (20) and the induction hypothesis} \\ &\leq p\left(n + \frac{(1-p)^2}{p}n\right) + (1 - p)pn && \text{by (18)} \\ &= pn + (1 - p)^2n + pn - p^2n = n \end{aligned}$$

Hence we have proved  $v_n \leq n$ . It follows that the expectation, starting in the top-left brown state of  $T^n$ , of the number of red *transitions* visited in  $T^n$  is at most  $n \cdot \frac{1}{n^2+1}$ . Thus the

expected number of visited red transitions in the whole MDP  $\mathcal{M}$  is at most  $\sum_{n=1}^{\infty} \frac{1}{n} \cdot n \cdot \frac{1}{n^2+1} \leq \frac{\pi^2}{6}$ . Hence there is  $k \in \mathbb{N}$  such that the expected number of red transitions visited in  $T^k, T^{k+1}, \dots$  is less than 1. It follows from the Markov inequality that the probability to visit at least one red transition in  $T^k, T^{k+1}, \dots$  is less than 1. Hence the probability to visit at least one red transition in  $\mathcal{M}$  is less than 1.

The strategy  $\sigma$  from above can be improved to obtain an  $\varepsilon$ -optimal strategy  $\sigma_\varepsilon$  for Büchi from  $s_0$ , i.e.,  $\mathcal{P}_{\mathcal{M}, s_0, \sigma_\varepsilon}(\varphi) \geq 1 - \varepsilon$ . We obtain  $\sigma_\varepsilon$  by modifying the described strategy  $\sigma$  such that, in the first  $k$  trees for some  $k \in \mathbb{N}$ , the upward transitions emanating from the yellow states are taken. By choosing a large but finite  $k$ , the risk of taking a red transition can be made arbitrarily small, while the probability of visiting infinitely many green states remains 1.  $\blacktriangleleft$

Finally, we are ready to prove our main claim, Theorem 3.

**Proof of Theorem 3.** We describe how to modify the MDP  $\mathcal{M}$  from Lemma 9 to obtain an MDP  $\mathcal{M}_2$  with the claimed properties. First eliminate the red transitions in  $\mathcal{M}$  and change the objective to the normal Büchi objective. This can be done by redirecting all red transitions to an infinite (losing) chain of non-green states. Denote the resulting MDP by  $\mathcal{M}_1$ . For a state  $s$ , define its *depth*  $d(s)$  as the length of the *longest* path from the start state  $s_0$  to  $s$ . In  $\mathcal{M}_1$  each state has finite depth (this property does not follow from acyclicity alone). Now obtain  $\mathcal{M}_2$  from  $\mathcal{M}_1$  by replacing every transition that leads from a state  $s_1$  to a state  $s_2$  with  $d(s_1) + 1 < d(s_2)$  by a chain (of non-green states) of length  $d(s_2) - d(s_1)$ . In this way, in  $\mathcal{M}_2$ , for every state  $s$ , all paths from  $s_0$  to  $s$  have the same length  $d(s)$ . Thus, instrumenting  $\mathcal{M}_2$  with a step-counter would lead to an MDP isomorphic to  $\mathcal{M}_2$ . It follows that every Markov strategy for  $\mathcal{M}_2$  could be replaced by an MR-strategy that achieves  $\text{Büchi}(F)$  with the same probability. Observe that a MR-strategy for  $\mathcal{M}_2$  directly translates to an MR-strategy for  $\mathcal{M}$  that achieves the same probability. Hence, item 1 follows, as the existence of a Markov-, and hence MR-strategy that achieves positive probability would contradict Lemma 9.

Item 2 is shown by modifying the strategies  $\sigma_\varepsilon$  from item 2 of Lemma 9 in the natural way.  $\blacktriangleleft$

A slight modification of the example above yields a lower bound on the memory requirements for the almost-sure parity objective. Recall that the parity objective is defined on systems whose states are labeled by a finite set of colors  $C \stackrel{\text{def}}{=} \{1, 2, \dots, \max\} \subseteq \mathbb{N}$ , where a run is in  $\text{Parity}(C)$  iff the highest color that is seen infinitely often in the run is even.

► **Corollary 4.** *There exist an acyclic MDP  $\mathcal{M}'$  with colors  $\{1, 2, 3\}$  and a state  $s_0$  such that*

1. *for every Markov strategy  $\sigma$ , we have  $\mathcal{P}_{\mathcal{M}', s_0, \sigma}(\text{Parity}(\{1, 2, 3\})) = 0$ , and*
2. *there exists a deterministic 1-bit strategy  $\sigma'$  such that  $\mathcal{P}_{\mathcal{M}', s_0, \sigma'}(\text{Parity}(\{1, 2, 3\})) = 1$ .*

**Proof.** We obtain  $\mathcal{M}'$  by modifying the MDP  $\mathcal{M}$  from Theorem 3 as follows. Label all green states in  $F$  by color 2 and the rest by color 1. Then modify each red transition to go to its target via a fresh state labeled by color 3. Clearly  $\mathcal{M}'$  is still acyclic and labeled by colors  $\{1, 2, 3\}$ .

From the proof of Theorem 3 (1), under every Markov strategy in  $\mathcal{M}$  a.s. seeing infinitely many green states (in  $F$ ) implies seeing infinitely many red transitions. So in  $\mathcal{M}'$  every Markov strategy  $\sigma$  a.s. either sees color 2 only finitely often or color 3 infinitely often, thus  $\mathcal{P}_{\mathcal{M}', s_0, \sigma}(\text{Parity}(\{1, 2, 3\})) = 0$ .

From the proof of [Theorem 3](#) (2), there is a deterministic 1-bit strategy  $\sigma$  in  $\mathcal{M}$  that attains probability  $\geq 1/2$  for  $\text{Büchi}(F)$  without taking any red transition and otherwise a.s. takes a red transition. This property of  $\sigma$  holds not only when starting from  $s_0$  but from every other state as well. We obtain  $\sigma'$  in  $\mathcal{M}'$  by continuing to play  $\sigma$  even after red transitions have been taken. Under  $\sigma'$  the probability of going through infinitely many red transitions (and seeing color 3) is  $\leq (1/2)^\infty = 0$ , and the probability of seeing infinitely many states in  $F$  (with color 2) is 1. Thus  $\mathcal{P}_{\mathcal{M}', s_0, \sigma'}(\text{Parity}(\{1, 2, 3\})) = 1$ .  $\blacktriangleleft$

## B The Upper Bound: Full Details

► **Theorem 5.** *For every acyclic countable MDP  $\mathcal{M}$ , finite set of initial states  $I$ , set of states  $F$  and  $\varepsilon > 0$ , there exists a deterministic 1-bit strategy for  $\text{Büchi}(F)$  that is  $\varepsilon$ -optimal from every  $s \in I$ .*

**Proof.** Let  $\mathcal{M} = (S, S_\square, S_\circ, \longrightarrow, P)$  be an acyclic MDP,  $I \subseteq S$  a finite set of initial states and  $F \subseteq S$  a set of goal states and  $\varphi \stackrel{\text{def}}{=} \text{Büchi}(F)$  denote the Büchi objective w.r.t.  $F$ . We prove the claim for finitely branching  $\mathcal{M}$  first and transfer the result to general MDPs at the end.

For every  $\varepsilon > 0$  and every  $s \in I$  there exists an  $\varepsilon$ -optimal strategy  $\sigma_s$  such that

$$\mathcal{P}_{\mathcal{M}, s, \sigma_s}(\varphi) \geq \text{val}_{\mathcal{M}, \varphi}(s) - \varepsilon. \quad (21)$$

However, the strategies  $\sigma_s$  might differ from each other and might use randomization and a large (or even infinite) amount of memory. We will construct a single deterministic strategy  $\sigma'$  that uses only 1 bit of memory such that  $\forall s \in I \mathcal{P}_{\mathcal{M}, s, \sigma'}(\varphi) \geq \text{val}_{\mathcal{M}, \varphi}(s) - 2\varepsilon$ . This proves the claim as  $\varepsilon$  can be chosen arbitrarily small.

In order to construct  $\sigma'$ , we first observe the behavior of the finitely many  $\sigma_s$  for  $s \in I$  on an infinite, increasing sequence of finite subsets of  $S$ . Based on this, we define a second stronger objective  $\varphi'$  with

$$\varphi' \subseteq \varphi, \quad (22)$$

and show that all  $\sigma_s$  attain at least  $\text{val}_{\mathcal{M}, \varphi}(s) - 2\varepsilon$  w.r.t.  $\varphi'$ , i.e.,

$$\forall s \in I \mathcal{P}_{\mathcal{M}, s, \sigma_s}(\varphi') \geq \text{val}_{\mathcal{M}, \varphi}(s) - 2\varepsilon. \quad (23)$$

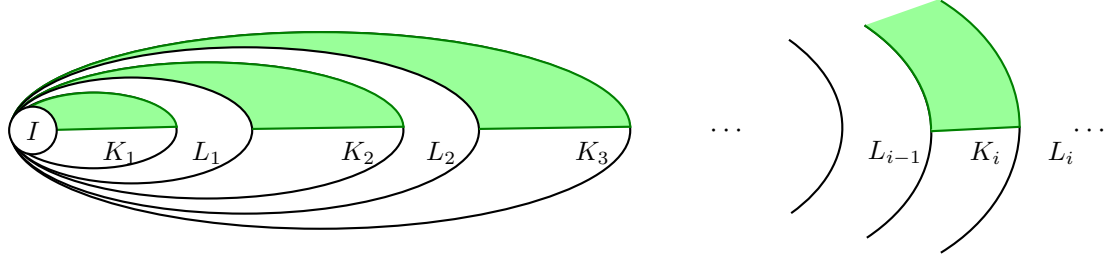
We construct  $\sigma'$  as a deterministic 1-bit *optimal* strategy w.r.t.  $\varphi'$  from all  $s \in I$  and obtain

$$\begin{aligned} \mathcal{P}_{\mathcal{M}, s, \sigma'}(\varphi) &\geq \mathcal{P}_{\mathcal{M}, s, \sigma'}(\varphi') && \text{by (22)} \\ &\geq \mathcal{P}_{\mathcal{M}, s, \sigma_s}(\varphi') && \text{by optimality of } \sigma' \text{ for } \varphi' \\ &\geq \text{val}_{\mathcal{M}, \varphi}(s) - 2\varepsilon && \text{by (23).} \end{aligned}$$

**Behavior of  $\sigma$ , objective  $\varphi'$  and properties (22) and (23).** We start with some notation. Let  $\text{bubble}_k(X)$  be the set of states that can be reached from some state in the set  $X$  within at most  $k$  steps. Since  $\mathcal{M}$  is finitely branching,  $\text{bubble}_k(X)$  is finite if  $X$  is finite. Let  $\text{Goal}^{\leq k}(X) \stackrel{\text{def}}{=} \{\rho \in S^\omega \mid \exists t \leq k. X(\rho(t)) = 1\}$  and  $\text{Goal}^{\geq k}(X) \stackrel{\text{def}}{=} \{\rho \in S^\omega \mid \exists t \geq k. X(\rho(t)) = 1\}$  denote the property of visiting the set  $X$  (at least once) within at most (resp. at least)  $k$  steps. Moreover, let  $\varepsilon_i \stackrel{\text{def}}{=} \varepsilon \cdot 2^{-(i+1)}$ .

The following lemma depends on the assumption that  $\mathcal{M}$  is acyclic.

► **Lemma 10.** *Let  $X \subseteq S$  be a finite set of states and  $\varepsilon' > 0$ .*



■ **Figure 5** To show the bubble construction. The green region in  $K_1$  is  $F_1$ , and for all  $i \geq 2$ , the green region in  $K_i \setminus L_{i-1}$  is  $F_i$ .

1. There is  $k \in \mathbb{N}$  such that  $\forall_{s \in I} \mathcal{P}_{\mathcal{M}, s, \sigma_s}(\varphi \cap \overline{\text{Goal}^{\leq k}(F \setminus X)}) \leq \varepsilon'$ .
2. There is  $l \in \mathbb{N}$  such that  $\forall_{s \in I} \mathcal{P}_{\mathcal{M}, s, \sigma_s}(\text{Goal}^{\geq l}(X)) \leq \varepsilon'$ .

**Proof.** It suffices to show the properties for a single  $s, \sigma_s$  since one can take the maximal  $k, l$  over the finitely many  $s \in I$ . By acyclicity of  $\mathcal{M}$ , it holds that  $\varphi \subseteq \text{Goal}(F \setminus X) = \bigcup_{k \in \mathbb{N}} \text{Goal}^{\leq k}(F \setminus X)$  and therefore that  $\varphi \cap \bigcap_{k \in \mathbb{N}} \overline{\text{Goal}^{\leq k}(F \setminus X)} = \emptyset$ . It follows from the continuity of measures that  $\lim_{k \rightarrow \infty} \mathcal{P}_{\mathcal{M}, s, \sigma_s}(\varphi \cap \overline{\text{Goal}^{\leq k}(F \setminus X)}) = 0$ .

Item 2 follows directly from the fact that  $\mathcal{M}$  is acyclic. ◀

By Lemma 10(1) there is a  $k_1$  such that for  $K_1 \stackrel{\text{def}}{=} \text{bubble}_{k_1}(I)$  and  $F_1 \stackrel{\text{def}}{=} F \cap K_1$  we have  $\forall_{s \in I} \mathcal{P}_{\mathcal{M}, s, \sigma_s}(\varphi \cap \overline{K_1^* F_1 S^\omega}) \leq \varepsilon_1$ . We define the pattern

$$R_1 \stackrel{\text{def}}{=} (K_1 \setminus F_1)^* F_1$$

and obtain  $\forall_{s \in I} \mathcal{P}_{\mathcal{M}, s, \sigma_s}(\varphi \cap \overline{R_1 S^\omega}) \leq \varepsilon_1$ . By Lemma 10(2) there is an  $l_1 > k_1$  such that  $\forall_{s \in I} \mathcal{P}_{\mathcal{M}, s, \sigma_s}(\text{Goal}^{\geq l_1}(K_1)) \leq \varepsilon_1$ . Define  $L_1 \stackrel{\text{def}}{=} \text{bubble}_{l_1}(I)$ . By Lemma 10(1) there is a  $k_2 > l_1$  such that for  $K_2 \stackrel{\text{def}}{=} \text{bubble}_{k_2}(I)$  and  $F_2 \stackrel{\text{def}}{=} F \cap K_2 \setminus L_1$  we have  $\forall_{s \in I} \mathcal{P}_{\mathcal{M}, s, \sigma_s}(\varphi \cap \overline{K_2^* F_2 S^\omega}) \leq \varepsilon_2$ . We define the pattern

$$R_2 \stackrel{\text{def}}{=} (K_2 \setminus F_2)^* F_2$$

and obtain  $\forall_{s \in I} \mathcal{P}_{\mathcal{M}, s, \sigma_s}(\varphi \cap \overline{R_2 S^\omega}) \leq \varepsilon_2$  and, via a union bound,  $\forall_{s \in I} \mathcal{P}_{\mathcal{M}, s, \sigma_s}(\varphi \cap \overline{R_2(S \setminus K_1)^\omega}) \leq \varepsilon_1 + \varepsilon_2$ . By another union bound it follows that  $\forall_{s \in I} \mathcal{P}_{\mathcal{M}, s, \sigma_s}(\varphi \cap \overline{R_1 R_2(S \setminus K_1)^\omega}) \leq 2\varepsilon_1 + \varepsilon_2$ .

Proceed inductively for  $i = 2, 3, \dots$  as follows (see Figure 5 for an illustration). By Lemma 10(2) there is an  $l_i > k_i$  such that  $\forall_{s \in I} \mathcal{P}_{\mathcal{M}, s, \sigma_s}(\text{Goal}^{\geq l_i}(K_i)) \leq \varepsilon_i$ . Define  $L_i \stackrel{\text{def}}{=} \text{bubble}_{l_i}(I)$ . By Lemma 10(1) there is  $k_{i+1} > l_i$  such that for  $K_{i+1} \stackrel{\text{def}}{=} \text{bubble}_{k_{i+1}}(I)$  and  $F_{i+1} \stackrel{\text{def}}{=} F \cap K_{i+1} \setminus L_i$  we have  $\forall_{s \in I} \mathcal{P}_{\mathcal{M}, s, \sigma_s}(\varphi \cap \overline{(K_{i+1} \setminus F_{i+1})^* F_{i+1} S^\omega}) \leq \varepsilon_{i+1}$ . By a union bound,  $\forall_{s \in I} \mathcal{P}_{\mathcal{M}, s, \sigma_s}(\varphi \cap \overline{(K_{i+1} \setminus F_{i+1})^* F_{i+1}(S \setminus K_i)^\omega}) \leq \varepsilon_i + \varepsilon_{i+1}$ . By an induction hypothesis we have  $\forall_{s \in I} \mathcal{P}_{\mathcal{M}, s, \sigma_s}(\varphi \cap \overline{R_1 R_2 \dots R_i(S \setminus K_{i-1})^\omega}) \leq 2\varepsilon_1 + \dots + 2\varepsilon_{i-1} + \varepsilon_i$ . We define the pattern

$$R_{i+1} \stackrel{\text{def}}{=} (K_{i+1} \setminus (F_{i+1} \cup K_{i-1}))^* F_{i+1}.$$

Using that  $(K_{i+1} \setminus F_{i+1})^* F_{i+1}(S \setminus K_i)^\omega \cap R_1 R_2 \dots R_i(S \setminus K_{i-1})^\omega \subseteq R_1 R_2 \dots R_{i+1}(S \setminus K_i)^\omega$ , we get

$$\forall_{s \in I} \mathcal{P}_{\mathcal{M}, s, \sigma_s}(\varphi \cap \overline{R_1 R_2 \dots R_{i+1}(S \setminus K_i)^\omega}) \leq 2\varepsilon_1 + \dots + 2\varepsilon_i + \varepsilon_{i+1} \leq \varepsilon. \quad (24)$$

We now define the Borel objectives  $R_{\leq i} \stackrel{\text{def}}{=} R_1 R_2 \dots R_i S^\omega$  and  $\varphi' \stackrel{\text{def}}{=} \bigcap_{i \in \mathbb{N}} R_{\leq i}$ . Since  $F_i \cap F_k = \emptyset$  for  $i \neq k$  and  $\varphi'$  implies a visit to the set  $F_i$  for all  $i \in \mathbb{N}$ , we have  $\varphi' \subseteq \varphi$  and



obtain (22). Moreover,  $R_{\leq 1} \supseteq R_{\leq 2} \supseteq R_{\leq 3} \dots$  is an infinite decreasing sequence of Borel objectives. For every  $s \in I$  we have

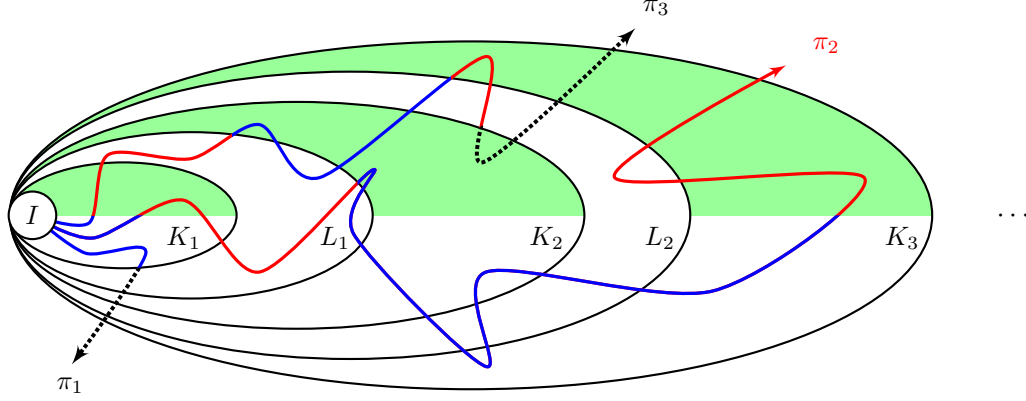
$$\begin{aligned}
\mathcal{P}_{\mathcal{M},s,\sigma_s}(\varphi') &= \mathcal{P}_{\mathcal{M},s,\sigma_s}(\cap_{i=1}^{\infty} R_{\leq i}) && \text{by def. of } \varphi' \\
&= \lim_{i \rightarrow \infty} \mathcal{P}_{\mathcal{M},s,\sigma_s}(R_{\leq i}) && \text{by cont. of measures} \\
&= \lim_{i \rightarrow \infty} 1 - \mathcal{P}_{\mathcal{M},s,\sigma_s}(\overline{R_{\leq i}}) && \text{by duality} \\
&= \lim_{i \rightarrow \infty} 1 - (\mathcal{P}_{\mathcal{M},s,\sigma_s}(\overline{R_{\leq i}} \cap \varphi) + \mathcal{P}_{\mathcal{M},s,\sigma_s}(\overline{R_{\leq i}} \cap \overline{\varphi})) && \text{case split} \\
&\geq \lim_{i \rightarrow \infty} 1 - (\varepsilon + \mathcal{P}_{\mathcal{M},s,\sigma_s}(\overline{R_{\leq i}} \cap \overline{\varphi})) && \text{by (24)} \\
&\geq \lim_{i \rightarrow \infty} 1 - (\varepsilon + \mathcal{P}_{\mathcal{M},s,\sigma_s}(\overline{\varphi'} \cap \overline{\varphi})) && \text{since } \varphi' \subseteq R_{\leq i} \\
&= 1 - (\varepsilon + 1 - \mathcal{P}_{\mathcal{M},s,\sigma_s}(\varphi' \cup \varphi)) && \text{by duality} \\
&= \mathcal{P}_{\mathcal{M},s,\sigma_s}(\varphi) - \varepsilon && \text{by (22)} \\
&\geq \text{val}_{\mathcal{M},\varphi}(s) - 2\varepsilon && \text{by (21)}
\end{aligned}$$

Thus we obtain property (23).

**Definition of the 1-bit strategy  $\sigma'$ .** We now define our deterministic 1-bit strategy  $\sigma'$  that is optimal for objective  $\varphi'$  from every  $s \in I$ . First we define certain “suffix” objectives of  $\varphi'$ . Recall that  $R_i = (K_i \setminus (F_i \cup K_{i-2}))^* F_i$ . Let  $R_{i,j} \stackrel{\text{def}}{=} R_i R_{i+1} \dots R_j S^\omega$  and  $R_{\geq i} \stackrel{\text{def}}{=} \bigcap_{j \geq i} R_{i,j}$ . In particular, this means that  $\varphi' = R_{\geq 1}$ . Every run  $w$  from some state  $s \in I$  that satisfies  $\varphi'$  can be split into parts before and after the first visit to set  $F_i$ , i.e.,  $w = w_1 s' w_2$  where  $w_1 s' \in R_{\leq i}$ ,  $s' \in F_i$  and  $s' w_2 \in R_{\geq i+1}$ . (Note also that  $w_2$  cannot visit any states in  $K_{i-1}$ .) Thus it will be useful to consider the objectives  $R_{\geq i+1}$  for runs that start in states  $s' \in F_i$ . For every state  $s' \in F_i$  we consider its value w.r.t. the objective  $R_{\geq i+1}$ , i.e.,  $\text{val}_{\mathcal{M},R_{\geq i+1}}(s') \stackrel{\text{def}}{=} \sup_{\hat{\sigma}} \mathcal{P}_{\mathcal{M},s',\hat{\sigma}}(R_{\geq i+1})$ .

For every  $i \geq 1$  we consider the finite subspace  $K_i \setminus K_{i-2}$ . In particular, it contains the sets  $F_{i-1}$  and  $F_i$ . (For completeness let  $K_0 \stackrel{\text{def}}{=} F_0 \stackrel{\text{def}}{=} I$  and  $K_{-1} \stackrel{\text{def}}{=} \emptyset$ .) It is not enough to maximize the probability of reaching the set  $F_i$  in each  $K_i$  individually. One also needs to maximize the potential of visiting further sets  $F_{i+1}, F_{i+2}, \dots$  in the indefinite future. Thus we define the bounded total reward objective  $B_i$  for runs starting in  $F_{i-1}$  as follows. Runs that exit the subspace (either by leaving  $K_i$  or by visiting  $K_{i-2}$ ) before visiting  $F_i$  get reward 0. All other runs must visit  $F_i$  eventually (since  $\mathcal{M}$  is acyclic and the subspace is finite). When some run reaches the set  $F_i$  for the first time in some state  $s'$  then this run gets the reward of  $\text{val}_{\mathcal{M},R_{\geq i+1}}(s')$ . We can consider an induced finite MDP  $\hat{\mathcal{M}}$  with state space  $K_i \setminus K_{i-2}$ , plus a sink state (with reward 0) that is reached immediately after visiting any state in  $F_i$  and whenever one exits the set  $K_i \setminus K_{i-2}$ . In  $\hat{\mathcal{M}}$  one gets a reward of  $\text{val}_{\mathcal{M},R_{\geq i+1}}(s')$  for visiting  $s' \in F_i$  as above. By [15, Theorem 7.1.9], there exists a uniform optimal MD-strategy  $\sigma_i$  for this bounded total reward objective on the induced finite MDP  $\hat{\mathcal{M}}$ , which can be directly applied for objective  $B_i$  on the subspace  $K_i \setminus K_{i-2}$  in  $\mathcal{M}$ . (The strategy  $\sigma_i$  is not necessarily unique, but our results hold regardless of which of them is picked.)

We now define  $\sigma'$  by combining different MD-strategies  $\sigma_i$ , depending on the current state and on the value of the 1-bit memory. The intuition is that the strategy  $\sigma'$  has two modes: normal-mode and next-mode. In a state  $s' \in K_i \setminus K_{i-1}$ , if the memory is  $i \pmod{2}$  then the strategy is in normal-mode and plays towards reaching  $F_i$ . Otherwise, the strategy is in next-mode and plays towards reaching  $F_{i+1}$  (normally this happens because  $F_i$  has already been seen).



■ **Figure 6** Memory updates along runs  $\pi_1, \pi_2, \pi_3$ , drawn in blue while the memory-bit is one and in red while the bit is zero. Both  $\pi_1$  and  $\pi_3$  violate  $\varphi'$  and are drawn as dotted lines once they do.

Initially  $\sigma'$  starts in a state  $s \in I$  with the 1-bit memory set to 1. We define the behavior of  $\sigma'$  in a state  $s' \in K_i \setminus K_{i-1}$  for every  $i \geq 1$ .

- If the 1-bit memory is  $i \pmod{2}$  and  $s' \notin F_i$  then  $\sigma'$  plays like  $\sigma_i$ . (Intuitively, one plays towards  $F_i$ , since one has not yet visited it.)
- If the 1-bit memory is  $i \pmod{2}$  and  $s' \in F_i$  then the 1-bit memory is set to  $(i+1) \pmod{2}$ , and  $\sigma'$  plays like  $\sigma_{i+1}$ . (Intuitively, one records the fact that one has already seen  $F_i$  and then targets the next set  $F_{i+1}$ .)
- If the 1-bit memory is  $(i+1) \pmod{2}$  then  $\sigma'$  plays like  $\sigma_{i+1}$ . (Intuitively, one plays towards  $F_{i+1}$ , since one has already visited  $F_i$ .)

Observe that if a run according to  $\sigma'$  exits some set  $K_i$  (and thus enters  $K_{i+1} \setminus K_i$ ) with the bit still set to  $i \pmod{2}$  (normal-mode) then this run has not visited  $F_i$  and thus does not satisfy the objective  $\varphi'$ . (Or the same has happened earlier for some  $j < i$ , in which case also the objective  $\varphi'$  is violated.) An example is the run  $\pi_1$  in Figure 6.

However, if a run according to  $\sigma'$  exits some set  $K_i$  (and thus enters  $K_{i+1} \setminus K_i$ ) with the bit set to  $(i+1) \pmod{2}$  (thus  $\sigma_{i+1}$  in next-mode) then in the new set  $K_{i'} \setminus K_{i'-1}$  with  $i' = i+1$  the bit is set to  $i' \pmod{2}$  and  $\sigma'$  continues to play like  $\sigma_{i+1}$  in normal-mode. Even if this run returns (temporarily) to  $K_i$  (but not to  $K_{i-1}$ ) the strategy  $\sigma'$  continues to play like  $\sigma_{i+1}$  in next-mode. An example is the run  $\pi_2$  in Figure 6.

Finally, if a run returns to  $K_{i-1}$  after having visited  $F_i$  then it fails the objective  $\varphi'$ . An example is the run  $\pi_3$  in Figure 6.

**The 1-bit strategy  $\sigma'$  is optimal for  $\varphi'$  from every  $s \in I$ .** In the following let  $s \in I$  be an arbitrary initial state in  $I$ . For any run from  $s$ , let  $\text{firstin}(F_i)$  be the first state  $s'$  in  $F_i$  that is visited (if any). We define a bounded reward objective  $B'_i$  for runs starting at  $s$  as follows. Every run that does not satisfy the objective  $R_{\leq i}$  gets assigned reward 0. Otherwise, consider a run from  $s$  that satisfies  $R_{\leq i}$ . When this run reaches the set  $F_i$  for the first time in some state  $s'$  then this run gets a reward of  $\text{val}_{\mathcal{M}, R_{\geq i+1}}(s')$ . Note that this reward is  $\leq 1$ .

We show that for all  $i \in \mathbb{N}$

$$\text{val}_{\mathcal{M}, \varphi'}(s) = \text{val}_{\mathcal{M}, B'_i}(s) \quad (25)$$

Towards the  $\geq$  inequality, let  $\hat{\sigma}$  be an  $\hat{\varepsilon}$ -optimal strategy for  $B'_i$  from  $s$ . We define the strategy  $\hat{\sigma}'$  to play like  $\hat{\sigma}$  until a state  $s' \in F_i$  is reached and then to switch to some  $\hat{\varepsilon}$ -optimal strategy for objective  $R_{\geq i+1}$  from  $s'$ . Every run from  $s$  that satisfies  $\varphi'$  can be split into parts, before and after the first visit to the set  $F_i$ , i.e.,  $\varphi' = \{w_1 s' w_2 \mid w_1 s' \in R_{\leq i}, s' \in F_i, s' w_2 \in R_{\geq i+1}\}$ . Therefore we obtain that  $\mathcal{P}_{\mathcal{M}, s, \hat{\sigma}'}(\varphi') \geq \mathcal{E}_{\mathcal{M}, s, \hat{\sigma}}(B'_i) - \hat{\varepsilon} \geq \text{val}_{\mathcal{M}, B'_i}(s) - 2\hat{\varepsilon}$ . Since this holds for every  $\hat{\varepsilon} > 0$ , we obtain  $\text{val}_{\mathcal{M}, \varphi'}(s) \geq \text{val}_{\mathcal{M}, B'_i}(s)$ .

Towards the  $\leq$  inequality, let  $\hat{\sigma}$  be any strategy for  $\varphi'$  from  $s$ . We have  $\mathcal{P}_{\mathcal{M}, s, \hat{\sigma}}(\varphi') \leq \sum_{s' \in F_i} \mathcal{P}_{\mathcal{M}, s, \hat{\sigma}}(R_{\leq i} \cap \text{firstin}(F_i) = s') \cdot \text{val}_{\mathcal{M}, R_{\geq i+1}}(s') = \mathcal{E}_{\mathcal{M}, s, \hat{\sigma}}(B'_i)$ . Thus  $\text{val}_{\mathcal{M}, \varphi'}(s) \leq \text{val}_{\mathcal{M}, B'_i}(s)$ . Together we obtain (25).

For all  $i \in \mathbb{N}$  and every state  $s' \in F_i$  we show that

$$\text{val}_{\mathcal{M}, R_{\geq i+1}}(s') = \text{val}_{\mathcal{M}, B_{i+1}}(s') \quad (26)$$

Towards the  $\geq$  inequality, let  $\hat{\sigma}$  be an  $\hat{\varepsilon}$ -optimal strategy for  $B_{i+1}$  from  $s' \in F_i$ . We define the strategy  $\hat{\sigma}'$  to play like  $\hat{\sigma}$  until a state  $s'' \in F_{i+1}$  is reached and then to switch to some  $\hat{\varepsilon}$ -optimal strategy for objective  $R_{\geq i+2}$  from  $s''$ . We have that  $\mathcal{P}_{\mathcal{M}, s', \hat{\sigma}'}(R_{\geq i+1}) \geq \mathcal{E}_{\mathcal{M}, s', \hat{\sigma}}(B_{i+1}) - \hat{\varepsilon} \geq \text{val}_{\mathcal{M}, B_{i+1}}(s') - 2\hat{\varepsilon}$ . Since this holds for every  $\hat{\varepsilon} > 0$ , we obtain  $\text{val}_{\mathcal{M}, R_{\geq i+1}}(s') \geq \text{val}_{\mathcal{M}, B_{i+1}}(s')$ .

Towards the  $\leq$  inequality, let  $\hat{\sigma}$  be any strategy for  $R_{\geq i+1}$  from  $s' \in F_i$ . We have

$$\begin{aligned} \mathcal{P}_{\mathcal{M}, s', \hat{\sigma}}(R_{\geq i+1}) &\leq \sum_{s'' \in F_{i+1}} \mathcal{P}_{\mathcal{M}, s', \hat{\sigma}}(R_{i+1} S^\omega \cap \text{firstin}(F_{i+1}) = s'') \cdot \text{val}_{\mathcal{M}, R_{\geq i+2}}(s'') \\ &= \mathcal{E}_{\mathcal{M}, s', \hat{\sigma}}(B_{i+1}). \end{aligned}$$

Thus  $\text{val}_{\mathcal{M}, R_{\geq i+1}}(s') \leq \text{val}_{\mathcal{M}, B_{i+1}}(s')$ . Together we obtain (26).

We show, by induction on  $i$ , that  $\sigma'$  is optimal for  $B'_i$  for all  $i \in \mathbb{N}$  from start state  $s$ , i.e.,

$$\mathcal{E}_{\mathcal{M}, s, \sigma'}(B'_i) = \text{val}_{\mathcal{M}, B'_i}(s) \quad (27)$$

In the base case of  $i = 1$  we have that  $B'_1 = B_1$ . The strategy  $\sigma'$  plays  $\sigma_1$  until reaching  $F_1$ , which is optimal for objective  $B_1$  and thus optimal for  $B'_1$ . For the induction step we assume (IH) that  $\sigma'$  is optimal for  $B'_i$ .

$$\begin{aligned} \text{val}_{\mathcal{M}, B'_{i+1}}(s) &= \text{val}_{\mathcal{M}, B'_i}(s) && \text{by (25)} \\ &= \mathcal{E}_{\mathcal{M}, s, \sigma'}(B'_i) && \text{by (IH)} \\ &= \sum_{s' \in F_i} \mathcal{P}_{\mathcal{M}, s, \sigma'}(R_{\leq i} \cap \text{firstin}(F_i) = s') \cdot \text{val}_{\mathcal{M}, R_{\geq i+1}}(s') && \text{by def. of } B'_i \\ &= \sum_{s' \in F_i} \mathcal{P}_{\mathcal{M}, s, \sigma'}(R_{\leq i} \cap \text{firstin}(F_i) = s') \cdot \text{val}_{\mathcal{M}, B_{i+1}}(s') && \text{by (26)} \\ &= \sum_{s' \in F_i} \mathcal{P}_{\mathcal{M}, s, \sigma'}(R_{\leq i} \cap \text{firstin}(F_i) = s') \cdot \mathcal{E}_{\mathcal{M}, s', \sigma_{i+1}}(B_{i+1}) && \text{opt. of } \sigma_{i+1} \text{ for } B_{i+1} \\ &= \mathcal{E}_{\mathcal{M}, s, \sigma'}(B'_{i+1}) && \text{by def. of } \sigma' \text{ and } B'_{i+1} \end{aligned}$$

So  $\sigma'$  attains the value  $\text{val}_{\mathcal{M}, B'_{i+1}}(s)$  of the objective  $B'_{i+1}$  from  $s$  and is optimal. Thus (27).

Now we show that  $\sigma'$  performs well on the objectives  $R_{\leq i}$  for all  $i \in \mathbb{N}$ .

$$\mathcal{P}_{\mathcal{M}, s, \sigma'}(R_{\leq i}) \geq \text{val}_{\mathcal{M}, \varphi'}(s) \quad (28)$$

We have

$$\begin{aligned} \mathcal{P}_{\mathcal{M}, s, \sigma'}(R_{\leq i}) &\geq \mathcal{E}_{\mathcal{M}, s, \sigma'}(B'_i) \quad \text{since } B'_i \text{ gives rewards 0 for runs } \notin R_{\leq i} \text{ and } \leq 1 \text{ otherwise} \\ &= \text{val}_{\mathcal{M}, B'_i}(s) \quad \text{by (27)} \\ &= \text{val}_{\mathcal{M}, \varphi'}(s) \quad \text{by (25)} \end{aligned}$$

So we get (28). Now we are ready to prove the optimality of  $\sigma'$  for  $\varphi'$  from  $s$ .

$$\begin{aligned}
\mathcal{P}_{\mathcal{M},s,\sigma'}(\varphi') &= \mathcal{P}_{\mathcal{M},s,\sigma'}(\cap_{i \in \mathbb{N}} R_{\leq i}) && \text{by def. of } \varphi' \\
&= \lim_{i \rightarrow \infty} \mathcal{P}_{\mathcal{M},s,\sigma'}(R_{\leq i}) && \text{by continuity of measures from above} \\
&\geq \lim_{i \rightarrow \infty} \text{val}_{\mathcal{M},\varphi'}(s) && \text{by (28)} \\
&= \text{val}_{\mathcal{M},\varphi'}(s)
\end{aligned}$$

This concludes the proof that  $\sigma'$  is optimal for  $\varphi'$  and hence  $2\varepsilon$ -optimal for  $\varphi$  for every initial state  $s \in I$ .

**From finitely to infinitely branching MDPs.** Let  $\mathcal{M}$  be an infinitely branching acyclic MDP with a finite set of initial states  $I$  and  $\varepsilon > 0$ . We derive a finitely branching acyclic MDP  $\mathcal{M}'$  with sufficiently similar behavior. Every controlled state  $x$  with infinite branching  $x \rightarrow y_i$  for all  $i \in \mathbb{N}$  is replaced by a gadget  $x \rightarrow z_1, z_i \rightarrow z_{i+1}, z_i \rightarrow y_i$  for all  $i \in \mathbb{N}$  with fresh controlled states  $z_i$  (cf. Figure 1). Infinitely branching random states with  $x \xrightarrow{p_i} y_i$  for all  $i \in \mathbb{N}$  are replaced by a gadget  $x \xrightarrow{1} z_1, z_i \xrightarrow{1-p'_i} z_{i+1}, z_i \xrightarrow{p'_i} y_i$  for all  $i \in \mathbb{N}$ , with fresh random states  $z_i$  and suitably adjusted probabilities  $p'_i$  to ensure that the gadget is left at state  $y_i$  with probability  $p_i$ , i.e.,  $p'_i = p_i / (\prod_{j=1}^{i-1} (1 - p'_j))$ .

We apply the above result for finitely branching acyclic MDPs to  $\mathcal{M}'$  and obtain a 1-bit deterministic  $\varepsilon$ -optimal strategy  $\sigma'$  for Büchi from all states  $s \in I$ . We construct a 1-bit deterministic  $\varepsilon$ -optimal strategy  $\sigma''$  for  $\mathcal{M}$  as follows. Consider some state  $x$  that is infinitely branching in  $\mathcal{M}$  and its associated gadget in  $\mathcal{M}'$ . Whenever a run in  $\mathcal{M}'$  according to  $\sigma'$  reaches  $x$  with some memory value  $\alpha \in \{0, 1\}$  there exist values  $p_i$  for the probability that the gadget is left at state  $y_i$ . Let  $p \stackrel{\text{def}}{=} 1 - \sum_{i \in \mathbb{N}} p_i$  be the probability that the gadget is never left. (If  $x$  is controlled then only one  $p_i$  (or  $p$ ) is nonzero, since  $\sigma'$  is deterministic. If  $x$  is random then  $p = 0$ .) Since  $\sigma'$  is deterministic, the memory updates are deterministic, and thus there are values  $\alpha'_i \in \{0, 1\}$  such that whenever the gadget is left at state  $y_i$  the memory will be  $\alpha'_i$ . We now define the behavior of the 1-bit deterministic strategy  $\sigma''$  at state  $x$  with memory  $\alpha$  in  $\mathcal{M}$ .

If  $x$  is controlled and  $p \neq 1$  then  $\sigma''$  picks the successor state  $y_i$  where  $p_i = 1$  and sets the memory to  $\alpha'_i$ . If  $p = 1$  then any run according to  $\sigma'$  that enters the gadget does not satisfy the objective. Thus  $\sigma''$  performs at least as well in  $\mathcal{M}$  regardless of its choice, e.g., pick successor  $y_1$  and  $\alpha' = \alpha$ .

If  $x$  is random then  $p = 0$  and the successor is chosen according to the defined distribution (which is the same in  $\mathcal{M}$  and  $\mathcal{M}'$ ) and  $\sigma''$  can only update its memory. Whenever the successor  $y_i$  is chosen,  $\sigma''$  updates the memory to  $\alpha'_i$ .

In states that are not infinitely branching in  $\mathcal{M}$ ,  $\sigma''$  does exactly the same in  $\mathcal{M}$  as  $\sigma'$  in  $\mathcal{M}'$ .

Since the gadgets do not intersect  $F$ ,  $\sigma''$  performs at least as well in  $\mathcal{M}$  as  $\sigma'$  in  $\mathcal{M}'$  and is thus  $\varepsilon$ -optimal from every  $s \in I$ .  $\blacktriangleleft$

Now we show our upper bound on the strategy complexity of Büchi for general MDPs.

► **Theorem 6.** *For every countable MDP  $\mathcal{M}$ , finite set of initial states  $I$ , set of states  $F$  and  $\varepsilon > 0$ , there exists a deterministic 1-bit Markov strategy for Büchi( $F$ ) that is  $\varepsilon$ -optimal from every  $s \in I$ .*

**Proof.** Let  $\mathcal{M} = (S, S_\square, S_\circ, \rightarrow, P)$  be a countable MDP with a finite set of initial states  $I$  and  $F \subseteq S$  the set of goal states. We derive an acyclic MDP  $\mathcal{M}' = (S', S'_\square, S'_\circ, \rightarrow', P')$

by encoding a step-counter into the states. Let  $S' = S \times \mathbb{N}_0$ ,  $S'_\square = S_\square \times \mathbb{N}_0$ ,  $S'_\circ = S_\circ \times \mathbb{N}_0$ ,  $\longrightarrow' = \{((x, n), (y, n+1)) \mid (x, y) \in \longrightarrow\}$  and  $P'((x, n))((y, n+1)) = P(x)(y)$  for all  $n \in \mathbb{N}_0$ . Let  $I' := \{(s, 0) \mid s \in I\}$  be the finite set of initial states of  $\mathcal{M}'$  and  $F' = F \times \mathbb{N}_0$  the set of goal states.

For every  $\varepsilon > 0$ , by Theorem 5, there exists a 1-bit deterministic strategy  $\sigma'$  for  $\text{Büchi}(F)$  in  $\mathcal{M}'$  that is  $\varepsilon$ -optimal from every state  $(s, 0) \in I'$ .

We now define the deterministic 1-bit Markov strategy  $\sigma$  for  $\text{Büchi}(F)$  that is  $\varepsilon$ -optimal from every  $s \in I$  in  $\mathcal{M}$ . It uses a step-counter (initially 0) and 1 extra bit of memory. For any controlled state  $s'$ , step-counter value  $n$  and memory  $\alpha \in \{0, 1\}$ , consider the behavior of  $\sigma'$  at state  $(s', n)$  and memory  $\alpha$ . Let  $(s'', n+1)$  be the chosen successor state and  $\alpha' \in \{0, 1\}$  the new memory content. Then  $\sigma$  chooses the successor  $s''$ , updates the memory to  $\alpha'$  and increments the step-counter to  $n+1$ .  $\blacktriangleleft$